MUSE: Secure Inference Resilient to Malicious Clients



Ryan Lehmkuhl UC Berkeley Pratyush Mishra UC Berkeley Akshayaram Srinivasan
Tata Institute of Fundamental Research*

Raluca Ada Popa UC Berkeley

Abstract

The increasing adoption of machine learning inference in applications has led to a corresponding increase in concerns about the privacy guarantees offered by existing mechanisms for inference. Such concerns have motivated the construction of efficient *secure inference* protocols that allow parties to perform inference without revealing their sensitive information. Recently, there has been a proliferation of such proposals, rapidly improving efficiency. However, most of these protocols assume that the client is semi-honest, that is, the client does not deviate from the protocol; yet in practice, clients are many, have varying incentives, and can behave arbitrarily.

To demonstrate that a malicious client can completely break the security of semi-honest protocols, we first develop a new *model-extraction attack* against many state-of-the-art secure inference protocols. Our attack enables a malicious client to learn model weights with $22\times-312\times$ fewer queries than the best black-box model-extraction attack [Car+20] and scales to much deeper networks.

Motivated by the severity of our attack, we design and implement MUSE, an efficient two-party secure inference protocol resilient to *malicious clients*. MUSE introduces a novel cryptographic protocol for *conditional disclosure of secrets* to switch between authenticated additive secret shares and garbled circuit labels, and an improved *Beaver's triple generation* procedure which is $8 \times -12.5 \times$ faster than existing techniques.

These protocols allow MUSE to push a majority of its cryptographic overhead into a preprocessing phase: compared to the equivalent *semi-honest* protocol (which is close to state-of-the-art), MUSE's online phase is only $1.7 \times -2.2 \times$ slower and uses $1.4 \times$ more communication. Overall, MUSE is $13.4 \times -21 \times$ faster and uses $2 \times -3.6 \times$ less communication than existing secure inference protocols which defend against malicious clients.

1 Introduction

The past few years have seen increasing deployment of neural network inference in popular applications such as image classification [Liu+17b] and voice assistants [Bar18]. However, the use of inference in such applications raises privacy concerns: existing implementations either require clients to send potentially sensitive data to remote servers for classification, or require the model owner to store their proprietary

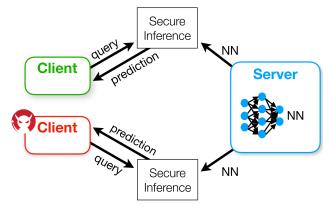


Figure 1: MUSE's system setup. Some of MUSE's clients may be malicious.

neural network model on the client's device. Both of these solutions are unsatisfactory: the former harms the privacy of the client, while the latter can harm a business model or reveal information about the training data and model weights.

To resolve this tension, the community has focused on constructing specialized protocols for *secure inference*, as we depict in Table 1. A secure inference protocol enables users and model owners to interact so that the user obtains the prediction result while ensuring that neither party learns any other information about the user input or the model weights. Many of these works implement these guarantees using secure two-party computation [Gil+16; Moh+17; Hes+17; Liu+17a; Bru+18; Cho+18; San+18; Juv+18; Lou+19; Bou+18; Rou+18; Bal+19; Ria+19; Dat+19; Mis+20; Rat+20]. However, as we can see from Table 1, all of these two-party works assume that both the client and the server follow the protocol rules, that is, they are *semi-honest*.

While it is common in the literature to assume a semihonest server, it is fundamentally less likely that *all* clients will behave correctly. The server is hosted at a single service provider, and existing cloud providers deploy competent intrusion detection systems, rigid access control, physical measures, and logging/tracking of the software installed [Goo17]. It is highly non-trivial to bypass these protections. Additionally, if a service provider is caught acting maliciously, the consequences may be dire due to public accountability.

In contrast, clients are many, run on a variety of setups under the control of users, users have various motives, and it suffices for only a single one of them to misbehave. The incentives for a client to cheat are high: service providers expend vast amounts of effort and money to accumulate data,

^{*}Work partially done while at UC Berkeley

			vulnerable to malicious clients	requires network ⁴ modification
	HE	CHET [Dat+19], CryptoDL [Hes+17], LoLa [Bru+18], TAPAS [San+18], Faster CryptoNets [Cho+18], FHE-DiNN [Bou+18]	•2	•
2PC		DeepSecure [Rou+18], [Bal+19] XONN [Ria+19]	•² •²	O
2	Mixed	SecureML [Moh+17] Gazelle [Juv+18], MiniONN [Liu+17a], CrypTFlow2 [Rat+20] DELPHI [Mis+20] MUSE	•¹ •¹ •¹ •¹	• • • •
3PC	Mixed	Chameleon [Ria+18] ABY ³ [Moh+18]	•1,3 O ³	0
	SS	SecureNN [Wag+19], Falcon [Wag+21], CrypTFlow [Kum+20]	O^3	0

Table 1: Related work on secure convolutional neural network (CNN) inference. See Section 7 for more details. This table compares specialized secure inference protocols, not generic frameworks for MPC. We compare against generic frameworks in Section 6. HE = Homomorphic Encryption, GC = Garbled Circuits, SS = Secret Sharing.

clean it, design model architectures, and to train the final model. If a client wishes to obtain a similar model, it would be much easier to steal the server's than to try and train an equivalent one; this makes *model-extraction attacks* attractive.

To illustrate the threat of a malicious client, in Section 2 we demonstrate a new model-extraction attack against semi-honest secure inference protocols whereby a malicious client can learn the server's entire model in a number of inference queries *linear* in the number of parameters in a network *regardless of its depth*. This attack outperforms the best model-extraction attacks for plaintext inference by 22×–312× [Jag+20; Car+20], and demonstrates that using semi-honest secure inference protocols can *significantly amplify* a malicious client's ability to steal a model.

A natural approach to defend against such an amplification is to leverage state-of-the-art generic secure computation tools providing *malicious security*. This approach guarantees that if *either* party acts maliciously, they will be caught and the protocol aborted, preserving privacy. However, such methods for achieving malicious security add a large overhead due to the use of heavy cryptographic primitives (e.g. zero-knowledge proofs [Gol+89] or cut-and-choose [Lin+15; Zhu+16]). In Section 6, we compare against such techniques.

To reduce this overhead, we propose MUSE, a secure inference protocol that works in the *client-malicious* threat model. In this model, the server is presumed to behave semi-honestly, but the client is allowed to deviate arbitrarily from the protocol description. As we will show in Section 6, working in this model enables MUSE to achieve much better performance than a fully malicious baseline.

Our contributions. To summarize, in this paper we make the following contributions:

• We devise a novel model-extraction attack against secure inference protocols that rely on additive secret sharing. This

- attack allows a malicious client to *perfectly extract all* the weights of a model with $22 \times -312 \times$ fewer queries than the state-of-the-art [Car+20]. The complexity of our attack depends *only* on the number of parameters, and not on other factors like the depth of the network.
- We present MUSE¹, an efficient two-party cryptographic inference protocol that is resilient to malicious clients. In designing MUSE, we develop a novel protocol for *conditional disclosure of secrets* to switch between authenticated additive secret shares and garbled circuit labels. Additionally, we formulate new client-malicious techniques for *triple generation* and *input authentication* in SPDZ-style MPC frameworks which improve performance by up to 12.5× and 37.8× respectively.
- Our implementation of MUSE is able to achieve an online phase that is only 1.7×-2.2× slower and uses 1.4× more communication than DELPHI [Mis+20], a recent protocol for *semi-honest* inference. When compared to fully-malicious secure inference protocols, MUSE is 13.4×-21× faster and uses 2×-3.5× less communication.

Remark 1.1. While MUSE's online phase is competitive with some of the best semi-honest protocols [Mis+20; Rat+20], the communication cost of preprocessing is up to $10 \times$ higher than in these semi-honest protocols. Hence, we view MUSE as a first step in constructing secure inference protocols that achieve client-malicious security, and anticipate that future works will rapidly lower this cost (the same has occurred for semi-honest secure inference protocols). MUSE already improves performance over current techniques for client-malicious inference by $13.4 \times -21 \times$ (see Section 6.4).

We now give a high-level overview of our techniques.

[•] Network modifications are optional ¹See Section 2.1 ² See Remark 2.1 ³Requires that two of the three parties act honestly ⁴Polynomial activations or binarized/discretized weights—may reduce network accuracy

¹MUSE is an acronym for Malicious-User Secure Inference

1.1 Our attack

What can a malicious client do? We start off by examining the power of malicious clients in secure inference protocols that rely on secret sharing. We noticed that many protocols of interest, such as [Moh+17; Juv+18; Liu+17a; Mis+20; Rat+20], have a similar structure, and we exploit this structure in our attack. The structure is as follows. These protocols "evaluate" the neural network in a layer-by-layer fashion, so that at the end of each layer, the client and the server both hold 2-out-of-2 secret shares of the output of that layer. At the end of the protocol, the server sends its share of the final output to the client, who uses it to reconstruct the final output.

Our attack relies on the following crucial observation: because the shares at the end of a layer are not authenticated, a malicious client can *additively malleate* them without detection. In more detail, let $\langle m \rangle_C$ be the client's share, and $\langle m \rangle_S$ be the server's share of a message $m \in \mathbb{F}$, so that $\langle m \rangle_C + \langle m \rangle_S = m$. Then, a malicious client can add an arbitrary shift r to a secret share to change the shared value from m to m+r. In Section 2, we show how one can leverage this malleability to learn the model weights.

1.2 Our protocol

We now explain how MUSE protects against malicious client attacks. We begin by describing our starting point: the semi-honest secure inference protocol DELPHI [Mis+20].

Starting point: DELPHI. We design MUSE by following the paradigm laid out in DELPHI [Mis+20]: since a convolutional neural network consists of alternating linear and non-linear layers, one should use subprotocols that are efficient for computing each type of layer, and then translate the output of one subprotocol to the input of the next. DELPHI instantiates these subprotocols by using additive secret sharing to evaluate linear layers, and garbled circuits to evaluate non-linear layers.

Attempt 1: Preventing malleability via MACs. The key insight in our attack in Sections 1.1 and 2 is that the client can malleate shares without detection. To prevent this, one can try to use standard techniques for authenticating the client's share via information-theoretic homomorphic message authentication codes (MACs). This technique is employed by the state-of-the-art protocols for malicious security [Kel+18; Che+20; Esc+20]. However, applying this technique directly to DELPHI runs into problems. For example, when switching between secret shares and garbled circuits, the server must ensure that the labels obtained by the client correspond to the authenticated secret share, and not to a different share. Doing this in a straightforward manner entails checking the share's MAC inside the garbled circuit, which is expensive. Furthermore, this check would need to be done in the online phase, which is undesirable.

Attempt 2: Separating authentication from computation. To remedy this, we make the following observation: garbled

circuits already achieve malicious security against garbled circuit evaluators (clients in our setting). This means that, if we had a specialized protocol that could output labels for the client's secret shares *only if* the corresponding MACs were valid, then we could compose this protocol with the garbled circuits to achieve an end-to-end client-malicious secure inference protocol.

In Section 5.1, we design exactly such a protocol for "conditional disclosure of secrets" (CDS). Unfortunately, executing our CDS protocol using existing frameworks for malicious MPC [Kel+18] proves to be extremely expensive. To address this, in Section 5.3 we devise a number of techniques to improve [Kel+18] in the client-malicious setting, and use the optimized framework to execute our CDS procedure. While the resulting protocol is much more efficient than checking MACs inside garbled circuits, it still imposes a significant cost on the online phase.

Our final protocol. To remedy this, our final insight in MUSE is that the secret shares and MACs that the client feeds into the CDS protocol *do not depend* on the client's input in the online phase. This allows us to move the execution of the CDS protocol *entirely to the preprocessing phase*, resulting in an online phase that is almost identical to that of DELPHI.

To summarize, in order to defend against malicious clients we first enforce authentication for the linear layers by using homomorphic MACs, then ensure that the client only receives garbled circuit labels corresponding to these authenticated shares via a novel CDS protocol, then develop new techniques for efficiently executing the CDS protocol, and finally move all these protocols to the preprocessing phase. For details, see Section 5.

2 Attacks on semi-honest inference protocols

We now describe how a malicious client can leverage the additive malleability of additive secret shares to learn the weights of a server's convolutional neural network in semi-honest secure inference protocols that rely on additive secret sharing. We begin in Section 2.1 by describing the kinds of protocols that are vulnerable to our attack. Then, in Section 2.2, we provide a detailed overview of our attack. Finally, in Section 2.3, we discuss the theoretical and empirical query complexity achieved by our attack.

2.1 Attack threat model

Our attack recovers the weights of neural networks consisting of alternating linear (that is, fully-connected or convolutional) and non-linear ReLU layers. Our attack works against semi-honest secure inference protocols that have the following properties:

 The protocol should evaluate the network iteratively by applying subprotocols for evaluating linear and non-linear layers.

²Our attack also supports networks with average pooling layers, as these are linear layers, but don't contain any weights that need to be recovered

- For each subprotocol, the input and output of the client and the server should be secret shares of the actual layer input and output, respectively.
- The client's final output should be the plaintext output of the final linear layer.

A number of two-party and multi-party secure inference protocols have these properties [Moh+17; Liu+17a; Ria+18; Juv+18; Cha+19; Mis+20; Rat+20].

Remark 2.1 (Other semi-honest protocols). Our attack does not affect semi-honest secure inference protocols based on fully-homomorphic encryption (FHE) [Gil+16; Cho+18; Bru+18; San+18; Dat+19] or garbled circuits (GC) [Rou+18; Bal+19; Ria+19]. However, this does not immunize these protocols against other kinds of malicious client attacks:

- FHE-based protocols use noise flooding [Gen09a] to hide the server's model. This technique is inherently semi-honest as it requires the pre-existing noise to be honestly bounded; if this does not hold, the noise term can reveal information about the server's model despite noise flooding.
- GC-based protocols use oblivious transfer (OT) [Rab81] to transfer labels for the client's input. However, if this OT is only semi-honest secure, a malicious client can attack it to learn both labels for the same input wire, which breaks the privacy guarantees of the garbled circuit, and leads to a leak of the server's model.

2.2 Attack strategy

Notation. Let NN be an ℓ -layer network convolutional neural network that classifies an image into one of m classes. That is, NN consists of ℓ matrices M_1, \ldots, M_ℓ so that $\mathsf{NN}(x) = M_\ell(\mathsf{ReLU}(\ldots M_2(\mathsf{ReLU}(M_1(x)))))$ where $M_\ell \in \mathbb{R}^{m \times t}$ and the image of M_ℓ is \mathbb{R}^m . We denote by $\mathsf{NN}_i(x)$ the partial evaluation of NN up to the i-th linear layer. That is, $\mathsf{NN}_i(x) := M_i(\mathsf{ReLU}(\ldots M_2(\mathsf{ReLU}(M_1(x)))))$. Below we denote by \mathbf{e}_j the j-th unit vector (the vector whose j-th entry is 1, and other entries are 0). Finally, for simplicity of exposition, we assume that biases are zero, 3 and that the network contains only fully-connected layers; for details on how to recover convolutional layers, see Remark 2.2 and Appendix A.

Prelude. Our attack proceeds in a bottom-up fashion: the client first recovers the parameters of the last linear layer $M_{\ell} \in \mathbb{R}^{m \times t}$, and then iteratively recovers previous layers. We describe the subroutine for recovering the last layer in Section 2.2.1, and then describe our subroutine for recovering intermediate layers in Section 2.2.2. In both subroutines, the client sets its initial input to the network be the all-zero vector.

2.2.1 Recovering the last layer

At a high level, to recover M_{ℓ} , the client proceeds columnby-column as follows: for each $j \in [t]$, the client provides as initial input the all-zero vector, and then honestly follows the secure inference protocol until the ℓ -th layer. At the ℓ -th layer, however, the client malleates its share of the input to M_{ℓ} so that it becomes \mathbf{e}_j . This means that result $M_{\ell} \cdot \mathbf{e}_j$ is the j-th column of M_{ℓ} . We illustrate this graphically for the first column below.

$$\begin{bmatrix}
\mathbf{e} \\
0 \\
0
\end{bmatrix}
\xrightarrow{\text{malleate}}
\xrightarrow{\mathbf{e}'_{\ell-1}}
\begin{bmatrix}
\mathbf{e} \\
1 \\
0
\end{bmatrix}
\xrightarrow{\mathbf{e}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
1 \\
0
\end{bmatrix}}
\xrightarrow{\mathbf{e}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
-0.1 \\
-1.1
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
-0.1 \\
0.2
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
0
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
-0.1 \\
-1.1
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
0
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
-0.1
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
0
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
-0.1
\end{bmatrix}}
\underbrace{\begin{bmatrix}
\mathbf{e} \\
0
\end{bmatrix}}
\underbrace{\begin{bmatrix}$$

2.2.2 Recovering intermediate layers

The foregoing algorithm works for recovering the last layer because the client can directly read off M_ℓ column-by-column by "solving" a linear system. However, this approach does not work as is for recovering the weights of intermediate linear layers, as we now demonstrate by considering the case of recovering the $\ell-1$ -th linear layer $M_{\ell-1}$. We then describe how to resolve the issues that arise. (The case of the remaining layers follows similarly).

Problem 1: Intervening ReLUs are lossy and non-linear. ReLUs between $M_{\ell-1}$ and M_{ℓ} disrupt the linearity of the system, preventing the use of linear system solvers.

Solution 1: Force ReLUs to behave linearly. To resolve this issue, we recall the fact that ReLU behaves like the identity function on inputs that are positive. We use malleability to exploit this property and *force* the remaining $M_{\ell}(\text{ReLU}(\cdot))$ computation to behave linearly, which means that we can once again solve a linear system to learn information about $M_{\ell-1}$.

In more detail, let $\langle \mathbf{y}_{\ell-2} \rangle_C$ be the client's share after applying $M_{\ell-1}$. The client malleates $\langle \mathbf{y}_{\ell-2} \rangle_C$ by setting it to $\langle \mathbf{y}'_{\ell-2} \rangle_C := \langle \mathbf{y}_{\ell-2} + \delta \rangle_C$, where δ is a constant vector whose elements are all greater than the magnitude of the largest element in $\mathbf{y}_{\ell-2}$. This forces all entries of $\mathbf{y}'_{\ell-2}$ to be positive, which means ReLU acts like the identity function. Then, after evaluating the ReLU and obtaining $\langle \mathbf{x}'_{\ell-1} \rangle_C$, the client "undoes" the malleation by subtracting δ . The following equation provides a graphical illustration of this process.

$$\overbrace{ \begin{bmatrix} -0.1 & 0.2 \\ -1.1 & 1.2 \end{bmatrix}}^{\textit{M}_{\ell-1}} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \overbrace{ \begin{bmatrix} -0.1 \\ -1.1 \end{bmatrix}}^{\textit{malleate}} \underbrace{ \begin{bmatrix} 9.9 \\ 8.9 \end{bmatrix}}^{\textit{malleate}} \underbrace{ \begin{bmatrix} x_{\ell-1} \\ 9.9 \\ 8.9 \end{bmatrix}}^{\textit{x}_{\ell-1}} \underbrace{ \underbrace{ \begin{bmatrix} x_{\ell-1} \\ 9.9 \\ 8.9 \end{bmatrix}}}_{\textit{unmalleate}} \underbrace{ \begin{bmatrix} -0.1 \\ -1.1 \end{bmatrix}}^{\textit{x}_{\ell-1}} \underbrace{ \begin{bmatrix} -0.1 \\ -1.1 \end{bmatrix}}^{\textit{malleate}} \underbrace{ \begin{bmatrix} -0.1 \\ -1.1 \end{bmatrix}}^{\textit{malleat$$

Problem 2: Underconstrained linear system. While the foregoing technique enables us to force the network to behave like a linear function, we have no guarantees that the resulting linear system is solvable. Indeed, neural networks necessarily map a high-dimensional feature to a low-dimensional classification, and so the resulting "linearized" neural network *must* be lossy. The following figure illustrates this graphically:

³One can handle a bias b in a linear layer L(x) = Mx + b by treating it as a simply another column in the modified matrix M' = M||b, so that the linear layer becomes $L(x) = M' \cdot (x||1)$.

⁴Note that since model weights are usually small (in the range [-1,1]), we can set δ to be a large value (say, \sim 10) to ensure that all entries of $\mathbf{y}'_{\ell-2}$ are positive.

$$\overbrace{ \begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}}^{M_2} \cdot \overbrace{ \begin{bmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \end{bmatrix}}^{M_1} = \overbrace{ \begin{bmatrix} a_1 + 3b_1 & a_2 + 3b_2 & a_3 + 3b_3 \\ 2a_1 + 4b_1 & 2a_2 + 4b_2 & 2a_3 + 4b_3 \end{bmatrix}}^{M_3}$$

Here, M_1 and M_2 are the first and last layers of the network, respectively. We have used the technique in Section 2.2.1 to recover M_2 , and now must recover M_1 . If we try to do this by querying M_3 , we get three (independent) equations for six variables, which is insufficient:

$$M_3\mathbf{e}_1 = \begin{bmatrix} 3a_1 + 6b_1 \\ 0 \\ 0 \end{bmatrix} \text{ and } M_3\mathbf{e}_2 = \begin{bmatrix} 0 \\ 3a_2 + 6b_2 \\ 0 \end{bmatrix} \text{ and } M_3\mathbf{e}_3 = \begin{bmatrix} 0 \\ 0 \\ 3a_3 + 6b_3 \end{bmatrix}$$

Solution 2: Masking variables. The issue is that M_3 does contain sufficient information to recover M_1 , but querying it naively loses that information. To resolve this, we use malleability again: the client uses the intervening ReLUs to "zero" out all but m entries of intermediate state, as follows:

$$M_1 \cdot \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} a_1 + a_2 \\ b_1 + b_2 \end{bmatrix} \xrightarrow{\text{malleate}} \begin{bmatrix} a_1 + a_2 + \delta \\ b_1 + b_2 - \delta \end{bmatrix} \xrightarrow{\text{ReLU +}} \begin{bmatrix} a_1 + a_2 \\ 0 \end{bmatrix}$$

Now, the client can obtain $M_2 \cdot [a_1 + a_2, 0]$, and can solve the resulting equations to learn a_1 and a_2 . It can then repeat this process with different queries and "masks" to learn all of M_1 .

For a detailed description of our algorithm, see Appendix A.

Remark 2.2 (recovering convolutional layers). To recover the kernel of a convolutional linear layer, we can reuse the foregoing ideas, but must change how we malleate the input to the target layer: we instead sample a random input and query the kernel via linearly independent columns of (the im2col transform of) this input. See Appendix A for details. Note that for simplicity of exposition, our description in Appendix A assumes that the number of channels and number of filters in each convolutional layer are both 1, and that the number of parameters in the kernel is less than the number of classes; these restrictions are easy to lift by adapting the masking techniques from above.

2.3 Efficiency and evaluation

Efficiency. Given a neural network where the *i*-th linear layer has dimension $m_i \times t_i$, and the number of classes is $m_\ell = m$, the foregoing algorithm learns the model parameters in just $\sum_{i=1}^{\ell} \lceil \frac{m_i}{m} \rceil \cdot t_i$ queries. Furthermore, the complexity of our attack depends *only on the number of parameters*, and not on other factors such as the depth. (This is not the case for other model-extraction attacks, which fail when extracting deep models that have few parameters.)

Evaluation. In Table 2, we compare our work to the state-of-the-art prior work on model extraction [Car+20], which does not rely on the existence of a secure inference protocol (and hence does not exploit properties of such protocols). Our experiments match the query complexity derived above.

network dimensions	# params	# queries		speedup
		us	[Car+20]	
FC-only networks				
784-128-1	100,480	100,480	2 ^{21.5}	29.5×
784-32-1	25,120	25,120	$2^{19.2}$	$24 \times$
10-10-10-1	210	210	2^{16}	312×
10-20-20-1	620	620	$2^{17.1}$	226.5×
40-20-10-10-1	1,110	1,110	$2^{21.5}$	205×
80-40-20-1	4,020	4,020	$2^{17.1}$	92×
80×5-40-20-1	29,620	29,620	_	n/a
1000-500-1	500,500	500,500	_	n/a
1000-500-10	505,000	50,500	_	n/a
2000×2-1000-100	6,100,000	61,000	_	n/a
1000×2-40×8-20-10	1,052,200	105,220	_	n/a
networks with convolu-	tions			
32 × (3,3)-P-100-10	158,088	15,809	_	n/a

Table 2: Query complexity of our attack vs. that of [Car+20]. The notation l-m-n-... indicates a series of fully-connected layers of dimension $l \times m$, $m \times n$, and so on, while $32 \times (3,3)$ indicates a convolutional layer consisting of $32 \times 3 \times 3$ filters, and P indicates a 2×2 average pooling layer.

3 Threat model and privacy goals

In our system, there are two parties: the client and the service provider (or server). The server holds a neural network model, and the client holds some data that it wants classified by the server's model. To achieve this goal, the two parties interact via a protocol for *secure inference*. This protocol takes as input the server's model and the client's data, and computes the classification so that neither party learns any information except this final classification. Below we clarify the security guarantees we aim for when designing our secure inference protocol MUSE.

3.1 Threat model

There are two standard notions of security for multiparty computation: security against semi-honest adversaries, and security against malicious adversaries. A semi-honest adversary follows the protocol perfectly but inspects messages it receives to learn information about other parties' inputs. A malicious adversary, on the other hand, may arbitrarily deviate from the protocol.

We design MUSE for a new threat model called "security against *malicious clients*" or *client-malicious security*. In this setting, either a malicious adversary corrupts the client, or a

⁵Note that any implementation of our attack will have to contend with errors due to limited floating point precision, but our experiments did not encounter such failures.

semi-honest adversary corrupts the server. ⁶

3.2 Privacy goals

MUSE's goal is to enable the client to learn at most the following information: the architecture of the neural network, and the result of the inference; all other information about the client's private inputs and the parameters of the server's neural network model should be hidden. Concretely, we aim to achieve a strong simulation-based definition of security as follows:

Definition 3.1. A protocol Π between a server and a client is said to securely compute a function f against a malicious client and semi-honest server if it satisfies the following properties:

- Correctness. For any server's input **y** and client's input **x**, the probability that at the end of the protocol, the client outputs $f(\mathbf{y}, \mathbf{x})$ is 1.
- Semi-Honest Server Security. For any server S that follows the protocol, there exists a simulator Sim_S such that for any input **y** of the server and **x** of the client, we have:

$$\mathsf{view}_S(\mathbf{y}, \mathbf{x}) \approx_c \mathsf{Sim}_S(\mathbf{y})$$

In other words, Sim_S is able to generate a view of the semi-honest server without knowing the client's private input.

 Malicious Client Security. For any malicious client C (that might deviate arbitrarily from the protocol specification), there exists a simulator Sim_C such that for any input y of the server, we have:

$$\mathsf{view}_C(\mathbf{y}) \approx_c \mathsf{Sim}_C^{f(\mathbf{y},\cdot)}$$

In other words, the Sim_C is able to generate the view of a malicious client with only access to an ideal functionality that accepts a client's input and outputs the result of the function f. This modeling is used in cryptographic literature to capture the cases where a malicious client may substitute its actual input with any other input of its choice.

Definition 3.2. We say that Π is a secure inference protocol against malicious clients and semi-honest servers if it securely computes $\mathsf{NN}(\cdot,\cdot)$ with the server input being \mathbf{M} and the client input being \mathbf{x} .

Like most prior work, MUSE does not hide information that is revealed by the result of the prediction. See Section 7.1 for a discussion of attacks that leverage this information, as well as potential mitigations.

4 Building blocks

MUSE uses the following cryptographic building blocks. Please see the full version for more formal definitions and proofs.

Garbling Scheme. A *garbling scheme* [Yao86; Bel+12] is a tuple of algorithms GS = (Garble, Eval) with the following syntax:

- Garble $(1^{\lambda}, C, \{\mathsf{lab}_{i,0}, \mathsf{lab}_{i,1}\}_{i \in [n]}) \to \tilde{C}$. On input the security parameter, a boolean circuit C (with n input wires) and a set of labels $\{\mathsf{lab}_{i,0}, \mathsf{lab}_{i,1}\}_{i \in [n]}$, Garble outputs a *garbled circuit* \tilde{C} . Here $\mathsf{lab}_{i,b}$ represents assigning the value $b \in \{0,1\}$ to the i-th input wire.
- Eval $(\tilde{C}, \{\mathsf{lab}_{i,x_i}\}_{i \in [n]}) \to y$. On input a garbled circuit \tilde{C} and labels $\{\mathsf{lab}_{i,x_i}\}_{i \in [n]}$ corresponding to an input $x \in \{0,1\}^n$, Eval outputs a string y = C(X).

We briefly describe here the key properties satisfied by garbling schemes. First, GS must be *correct*: the output of Eval must equal C(x). Second, it must be *private*: given \tilde{C} and $\{\mathsf{lab}_{i,x_i}\}$, the evaluator should not learn anything about C or x except the size of |C| (denoted by $1^{|C|}$) and the output C(x).

Leveled Fully Homomorphic public-key encryption.

A *leveled fully-homomorphic* encryption scheme HE = (KeyGen, Enc, Dec, Eval) [Reg09; Fan+12] is a public key encryption scheme that additionally supports homomorphically evaluating any depth-*D* arithmetic circuit on encrypted messages. Formally, HE satisfies the following syntax and properties:

- KeyGen(1^λ) → (pk,sk): On input a security parameter, KeyGen outputs a public key pk and a secret key sk.
- Enc(pk,m) $\rightarrow c$: On input the public key pk and a message m, the encryption algorithm Enc outputs a ciphertext c. We assume that the message space is \mathbb{Z}_p for some prime p.
- Dec(sk,c) → m: On input a secret key sk and a ciphertext c, the decryption algorithm Dec outputs a message m.
- Eval(pk, c₁, c₂, f) → c': On input a public key pk, ciphertexts c₁ and c₂ encrypting m₁ and m₂ respectively, and a depth-D arithmetic circuit f, Eval outputs a new ciphertext c'.

Besides the standard correctness and semantic security properties, we require HE to satisfy the following properties:

- Homomorphism. If $c_1 := \mathsf{Enc}(\mathsf{pk}, m_1), c_2 := \mathsf{Enc}(\mathsf{pk}, m_2),$ and $c := \mathsf{Eval}(\mathsf{pk}, c_1, c_2, f),$ then $\mathsf{Dec}(\mathsf{sk}, c) = f(m_1, m_2).$
- Function privacy. Given a ciphertext c, no attacker can tell what homomorphic operations led to c.

Additive secret sharing. Let p be a prime. A 2-of-2 additive secret sharing of $x \in \mathbb{Z}_p$ is a pair $(\langle x \rangle_1, \langle x \rangle_2) = (x - r, r) \in \mathbb{Z}_p^2$ for a random $r \in \mathbb{Z}_p$ such that $x = \langle x \rangle_1 + \langle x \rangle_2$. Additive secret sharing is perfectly hiding, i.e., given a share $\langle x \rangle_1$ or $\langle x \rangle_2$, the value x is perfectly hidden.

Message authentication codes. A *message authentication code (MAC)* is a tuple of algorithms MAC = (KeyGen, Tag, Verify) with the following syntax:

⁶One can generalize this threat model to *n* parties by considering two fixed subsets of parties: one of which can be corrupted by a malicious adversary, and the other which can be corrupted by a semi-honest adversary

- KeyGen $(1^{\lambda}) \rightarrow \alpha$: On input the security parameter, KeyGen outputs a MAC key α .
- Tag(α ,m) \rightarrow σ : On input a key α and message m, Tag outputs a tag σ and a secret state st.
- Verify(α, st, m, σ) → {0,1}: On input a key α, secret state st, message m and tag σ, Verify outputs 0 or 1.

We require MAC to satisfy the following properties:

- *Correctness*. For any message m, $\alpha \leftarrow \text{KeyGen}(1^{\lambda})$, and $(\sigma, \text{st}) \leftarrow \text{Tag}(\alpha, m)$, $\text{Verify}(\alpha, \text{st}, m, \sigma) = 1$.
- *One-time Security*. Given a valid message-tag pair, no adversary can forge a different, valid message-tag pair.

In this work, we will use the following construction of MACs:

- 1. The message space is \mathbb{Z}_p^n for some $p^n \ge 2^{\lambda}$.
- 2. KeyGen samples a uniform element $\alpha \leftarrow \mathbb{Z}_p$.
- 3. Tag(α , m) outputs $\sigma = \langle \alpha \cdot m \rangle_1$ and st = $\langle \alpha \cdot m \rangle_2$.
- 4. Verify(α , st, m, σ) checks if σ + st = $\alpha \cdot m$.

Beaver's multiplicative triples. A multiplication triple is a triple $(a,b,c) \in \mathbb{Z}_p^3$ such that ab = c. A triple generation procedure is a two-party protocol that outputs secret shares of a triple (a,b,c) to two parties.

Authenticated secret shares. For any prime p, an element $x \in \mathbb{Z}_p$, and a MAC key $\delta \in \mathbb{Z}_p$ an *authenticated share* of x is a tuple $(\varepsilon, [\![x]\!]_1, [\![x]\!]_2) := (\varepsilon, (\langle x \rangle_1, \langle \delta \cdot x \rangle_1), (\langle x \rangle_2, \langle \delta \cdot x \rangle_2))$. An authenticated share naturally supports local evaluation of addition and multiplication by public constants, as well as addition with another authenticated share. To multiply two authenticated shares, one needs to use multiplication triples. For simplicity of exposition, in the rest of the paper we omit ε , as it is merely used for bookkeeping when adding public constants.

Zero-knowledge proofs. Let \mathcal{R} be any NP relation. A *zero-knowledge proof* for \mathcal{R} is a protocol between a prover P and a verifier V that both have a common input x, where P tries to convince V that it "knows" a secret witness w such that $(x,w) \in \mathcal{R}$. At the end of the protocol, V should have learnt no additional information about w. We want our zero-knowledge proof system to satisfy the standard definitions of *completeness*, *soundness*, *proof of knowledge*, and *zero-knowledge*.

5 The MUSE protocol

In this section, we describe MUSE, our secure inference protocol that is secure against a malicious client and a semi-honest server. Like the DELPHI protocol (see Fig. 2) [Mis+20], MUSE's protocol consists of two phases: an offline preprocessing phase, and an online inference phase. The offline preprocessing phase is independent of the client's input (which regularly changes), but assumes that the server's model is static; if this model changes, then both parties have to rerun the preprocessing phase. After preprocessing, during the online inference phase, the client provides its input to our specialized secure two-party computation protocol, and eventually learns the inference result. Below, we expand on the

Preprocessing phase. During preprocessing, the client and the server pre-compute data for the online execution. This phase can be executed independently of the input values, i.e., DELPHI can run this phase before either party's input is known. Preprocessed data can only be used for a single inference.

- 1. *Linear correlations generator:* The client and server interact with a functionality that, for each $i \in [\ell]$, outputs to them secret shares of $M_i \mathbf{r}_i$, where \mathbf{r}_i is a random masking vector.
- 2. Preprocessing for ReLUs: The server constructs a garbled circuit \widetilde{C} for a circuit C computing ReLU. It sends \widetilde{C} to the client and then uses OT to send to the client the input wires corresponding to \mathbf{r}_{i+1} and $\mathbf{M}_i \cdot \mathbf{r}_i \mathbf{s}_i$.

Online phase. The online phase is divided into two stages:

- Preamble: On input x, the client sends x r₁ to the server.
 The server and the client now hold an additive secret sharing of x
- 2. Layer evaluation: Let \mathbf{x}_i be the result of evaluating the first (i-1) layers of the neural network on \mathbf{x} . At the beginning of the i-th layer, the client holds \mathbf{r}_i , and the server holds $\mathbf{x}_i \mathbf{r}_i$, which means that they possess secret shares of \mathbf{x}_i .
 - Linear layer: The server computes M_i · (x_i r_i), which
 means that the client and the server hold an additive secret
 sharing of M_ix_i.
 - ReLU layer: After the linear layer, the client and server hold secret shares of $M_i \mathbf{x}_i$. The server sends to the client the labels corresponding to its secrete share, and the client then evaluates the GC to obtain a secret share of the ReLU output.

Figure 2: High-level overview of the DELPHI protocol [Mis+20].

high level overview in Section 1.2 and provide a detailed description of both phases of our protocol.

Notation. The server holds a model \mathbf{M} consisting of ℓ linear layers $\mathbf{M}_1, \dots, \mathbf{M}_\ell$ and the client holds an input vector $\mathbf{x} \in \mathbb{Z}_p^n$. We use $\mathsf{NN}(\mathbf{M}, \mathbf{x})$ to denote the output of the neural network when the server's input is \mathbf{M} and the client's input is \mathbf{x} . We assume that the algorithm computing NN is public and is known to both the client and the server.

5.1 Preprocessing phase

In the preprocessing phase, the client and the server precompute data that can be used during the online execution. This phase is independent of the client's input values, and can be run before the client's input is known. However, this phase cannot be reused and has to be run once for each client input.

5.1.1 Intuition

As explained in Section 1.2, MUSE follows the approach of DELPHI, and uses different cryptographic primitives to produce preprocessed material for linear and non-linear layers. Below we describe these primitives at a high level.

Linear layers. Like in DELPHI, our goal is to produce shares of $M\mathbf{r}$ for a linear layer M. This enables us to efficiently compute linear layer operations in the online phase. Unlike

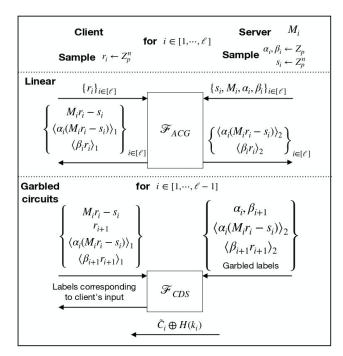


Figure 3: MUSE preprocessing phase.

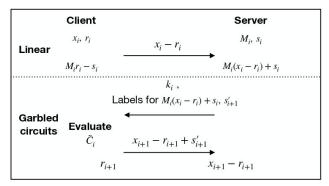


Figure 4: MUSE online phase.

DELPHI, we additionally need to prevent tampering by malicious clients. To this end, we extend the linear correlations generator (CG) used in DELPHI (see Fig. 2) to additionally support authentication. We formalize this via functionality \mathcal{F}_{ACG} for generating authenticated correlations. See Fig. 11 for a formal description.

To construct a protocol Π_{ACG} that realizes \mathcal{F}_{ACG} , we extend the techniques based on the leveled fully-homomorphic encryption used in DELPHI to additionally authenticate and secret share the relevant ciphertexts (see Fig. 5). We also require the client to provide zero-knowledge proofs that assert that their input ciphertexts are well-formed.

Non-linear layers. Like in DELPHI, we use garbled circuits to efficiently evaluate ReLUs. However, unlike DELPHI, we can no longer use oblivious transfer to send garbled labels to the client, because we have no way to check that the input to the oblivious transfer corresponds to the output from \mathcal{F}_{ACG} . Instead we introduce a functionality which conditionally out-

puts these labels if the inputs match the output of \mathcal{F}_{ACG} . We call this functionality *Conditional Disclosure of Secrets*, and denote it by \mathcal{F}_{CDS} .

To construct a protocol Π_{CDS} that realizes \mathcal{F}_{CDS} , we have two options: use 2PC protocols specialized for boolean computation, or 2PC protocols specialized for arithmetic computation. Indeed, because this operation fundamentally reasons about boolean values, it would seem reasonable to use a protocol like garbled circuits. However, checking validity of the client's input requires modular multiplications, which are extremely expensive when expressed as boolean circuits. Since even the simplest neural networks oftentimes have thousands of activations, the resulting communication and computation cost is unacceptable.

Instead, we implement this functionality via MPC for arithmetic circuits, as modular multiplication is cheap here. However, now the boolean operations are expensive. To overcome this, we take further advantage of the client-malicious setting to improve the MPC protocol we use to securely execute the arithmetic circuit, as we describe in Section 5.3.

By designing efficient protocols for \mathcal{F}_{ACG} and \mathcal{F}_{CDS} , MUSE achieves client-malicious security with an online phase design *identical to that of the semi-honest* DELPHI *protocol*. In our implementation, there are a few differences we detail in Remarks 5.2 and 5.3.

5.1.2 Protocol

We now present the full protocol for the preprocessing phase of MUSE (see Fig. 3 for a graphical overview).

- 1. For every $i \in [\ell]$, denote n_i, m_i as the input and output sizes of the i-th linear layer respectively. The client samples a random layer input mask $\mathbf{r}_i \leftarrow \mathbb{Z}_p^{n_i}$ and the server samples a random layer output mask $\mathbf{s}_i \leftarrow \mathbb{Z}_p^{m_i}$. Additionally, the server samples random MAC keys $\alpha_i, \beta_i \leftarrow \mathbb{Z}_p$.
- 2. **Authenticated correlations generator:** The client and server invoke functionality \mathcal{F}_{ACG} with the client input $\{\mathbf{r}_i\}_{i\in[\ell]}$, and with server input $\{\mathbf{s}_i,\mathbf{M}_i,\alpha_i,\beta_i\}_{i\in[\ell]}$. For each $i\in[\ell]$, the client obtains $\{\mathbf{M}_i(\mathbf{r}_i)-\mathbf{s}_i\}_{i\in[\ell]}$ along with $\{\langle \beta_i \cdot \mathbf{r}_i \rangle_1, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i)-\mathbf{s}_i \rangle)_1\}$ whereas the server receives $\{\langle \beta_i \cdot \mathbf{r}_i \rangle_2, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i)-\mathbf{s}_i \rangle)_2\}$. In Fig. 11 we describe the ideal functionality in more detail, and give a protocol for achieving it in Fig. 5.
- 3. For each $i \in [\ell]$, let $\mathsf{inp}_i := (\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i, \mathbf{r}_{i+1})$ denote the client's input to the i-th non-linear layer and $|\mathsf{inp}_i|$ denote its size in bits. The server chooses a set of random garbled circuit input labels $\{\mathsf{lab}_{i,k,0}^C, \mathsf{lab}_{i,k,1}^C\}_{k \in [|\mathsf{inp}_i|]}$.
- 4. Conditional disclosure of secrets: For each $i \in [\ell]$, the client and the server invoke functionality \mathcal{F}_{CDS} on the client's input inp_i and the MAC shares received from \mathcal{F}_{ACG} . If the client honestly inputs the correct shares, the functionality outputs the garbled input labels $\{lab_{i,k,inp_i}^C\}_{k \in [[inp_i]]}$ corresponding to inp_i to the client. In Fig. 12, we describe the ideal functionality in more detail, and give a protocol for securely computing this functionality in Fig. 6.

- 5. For each $i \in [\ell]$, the server chooses random labels $\{\mathsf{lab}_{i,k,0}^S,\mathsf{lab}_{i,k,1}^S\}_{k\in[|\mathsf{inp}_i|]}$ for its input to the i^{th} non-linear laver.7
- 6. Offline garbling: For each $i \in [\ell]$, the server garbles the circuit C_i (described in Fig. 7) using $\{\mathsf{lab}_{i,k,0}^C,\mathsf{lab}_{i,k,1}^C,\mathsf{lab}_{i,k,0}^S,\mathsf{lab}_{i,k,1}^S\}_{k\in[[\mathsf{inp}_i]]} \text{ as the input labels }$ to obtain the garbled circuit \widetilde{C}_i . It chooses a key $k_i \leftarrow$ $\{0,1\}^{\lambda}$ and sends $H(k_i) \oplus \widetilde{C}_i$ to the client where H is a random oracle.

5.2 Online phase

The online phase is divided into two stages: the *preamble* and the *layer evaluation*. (See Fig. 4 for a graphical overview.)

5.2.1 Preamble

The client sends $(\mathbf{x} - \mathbf{r}_1)$ to the server.

5.2.2 Layer evaluation

At the beginning of evaluating the *i*-th layer, the client holds \mathbf{r}_i and the server holds $\mathbf{x}_i - \mathbf{r}_i$ where \mathbf{x}_i is the vector obtained by evaluating the first (i-1) layers of the neural network on input x (with $x_1 = x$). This invariant will be maintained for each layer. We now describe the protocol for evaluating the i-th layer, which consists of linear functions and activation functions:

Linear layer. The server computes $\mathbf{M}_i(\mathbf{x}_i - \mathbf{r}_i) + \mathbf{s}_i$ which ensures that the client and server hold an additive secret share of $\mathbf{M}_{i}\mathbf{x}_{i}$.

Non-linear layer. After the linear layer, the server holds $\mathbf{M}_i(\mathbf{x}_i - \mathbf{r}_i) + \mathbf{s}_i$ and the client holds $\mathbf{M}_i \mathbf{r}_i - \mathbf{s}_i$. The parties evaluate the non-linear garbled circuit layer as follows:

- 1. The server chooses a random masking vector \mathbf{s}'_{i+1} and sends the labels from the set $\{\mathsf{lab}_{i,k,0}^S, \mathsf{lab}_{i,k,1}^S\}_{k \in [|\mathsf{inp}_i|]}$ corresponding to its input $\mathbf{M}_i(\mathbf{x}_i - \mathbf{r}_i) + \mathbf{s}_i$ and \mathbf{s}'_{i+1} to the client along with the key k_i .
- 2. The client uses k_i to unmask $H(k_i) \oplus \widetilde{C}_i$ to obtain \widetilde{C}_i . The client evaluates the garbled circuit C_i using its input labels obtained in the preprocessing phase and labels obtained from the server in the online phase. The client decodes the output labels and sends $\mathbf{x}_{i+1} - \mathbf{r}_{i+1} + \mathbf{s}'_{i+1}$ along with the hash of the output labels to the server.
- 3. The server checks if the hash computation is correct and recovers $\mathbf{x}_{i+1} - \mathbf{r}_{i+1}$.

Output phase. The server sends $\mathbf{s}'_{\ell+1}$ to the client and the client unmasks the output of the garbled circuit using this to learn the output of the inference y.

Theorem 5.1. Assuming the security of garbled circuits and the protocols for securely computing \mathcal{F}_{ACG} and \mathcal{F}_{CDS} , the protocol described above is a private inference protocol against malicious clients and semi-honest servers (see Definition 3.2) in the random oracle model.

Protocol Π_{ACG}

- 1. Both parties engage in a two-party computation protocol with security against malicious clients and semi-honest servers to generate (pk,sk) for HE. The client learns pk and sk whereas the server only learns pk.
- 2. The client sends $\{\mathsf{Enc}(\mathsf{pk},\mathbf{r}_i)\}_{i\in[\ell]}$ to the server along with a zero-knowledge proof of well-formedness of the ciphertext. The server verifies this proof before continuing.
- 3. For every $i \in [\ell]$,
 - (a) The homomorphically computes server $\mathsf{Enc}(\mathsf{pk}, \mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i), \ \mathsf{Enc}(\mathsf{pk}, \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i)),$ and Enc(pk, $\beta_i \cdot \mathbf{r}_i$).
 - (b) The server randomly samples $\langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_2$ and $\langle \beta_i \cdot \mathbf{r}_i \rangle_2$, homomorphically creates additive shares of the MAC values, and sends $(Enc(pk, M_i(r_i) (\mathbf{s}_i), \mathsf{Enc}(\mathsf{pk}, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i) \rangle_1), \mathsf{Enc}(\mathsf{pk}, \langle \beta_i \cdot \mathbf{r}_i \rangle_1))$ to the client.
 - (c) The client decrypts the above ciphertexts and obtains $(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i), \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i) \rangle_1$, and $\langle \beta_i \cdot \mathbf{r}_i \rangle_1$. The server holds $\langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i) \rangle_2$ and $\langle \beta_i \cdot \mathbf{r}_i \rangle_2$.

Figure 5: Our construction of an authenticated correlations generator (ACG).

Protocol Π_{CDS}

Denote the bit decomposition of $inp_i = (\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i, \mathbf{r}_{i+1})$ as

- The client and server input securely compute the following function f_{CDS} using a secure two-party computation protocol against malicious clients and semi-honest servers:
 - input consists $(\{b_k^i\}_{k \in [[\mathsf{inp}_i]]}, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i) \rangle_1, \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_1) \text{ and the server's input consists of } (\alpha_i, \beta_{i+1}, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i) \rangle_1) = (\alpha_i, \beta_{i+1}, \alpha_i(\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i))$ $\begin{array}{l} \mathbf{s}_i)\rangle_2, \langle \beta_{i+1} \cdot \mathbf{r}_{i+1}\rangle_2, \{\mathsf{lab}_{i,k,0}^C, \mathsf{lab}_{i,k,1}^C\}_{k \in [[\mathsf{inp}_i]]}) \\ \mathsf{some} \ i \in [\ell-1]. \end{array}$
 - (b) Denote $g_k = \mathsf{lab}_{i,k,0}^C (\mathsf{lab}_{i,k,0}^C \mathsf{lab}_{i,k,1}^C) \cdot b_k^i$. Additively secret share g_k into $(\langle g_k \rangle_1, \langle g_k \rangle_2)$.
 - (c) Reconstruct $\overline{\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i}$ and $\overline{\mathbf{r}_{i+1}}$ from $\{b_k^t\}_{k \in [|\mathsf{inp}_i|]}$.
 - (d) Sample random vectors $\mathbf{c}_1 \leftarrow \mathbb{Z}_p^{|\mathbf{r}_{i+1}|}, \mathbf{c}_2 \leftarrow \mathbb{Z}_p^{|\mathbf{r}_{i+1}|}$
 - (e) Denote $\rho = \alpha_{i}(\mathbf{c}_{1}^{T} \cdot \overline{\mathbf{M}_{i}}(\mathbf{r}_{i}) \mathbf{s}_{i}) \mathbf{c}_{1}^{T} \cdot (\langle \alpha_{i}(\mathbf{M}_{i}(\mathbf{r}_{i}) \mathbf{s}_{i}) \rangle_{1} + \langle \alpha_{i}(\mathbf{M}_{i}(\mathbf{r}_{i}) \mathbf{s}_{i}) \rangle_{2}) \text{ and } \sigma = \beta_{i+1}(\mathbf{c}_{2}^{T} \cdot \overline{\mathbf{r}_{i+1}}) \mathbf{c}_{2}^{T} \cdot (\langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_{1} + \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_{2}).$ (f) Output ρ , σ , $\{\langle g_{k} \rangle_{2}\}_{k \in [[\inf \rho_{i}]]}$ to the server and
 - $\{\langle g_k \rangle_1\}_{k \in [|\mathsf{inp}_i|]}$ to the client.
- 2. If either of them are non-zero, the server aborts the protocol. Else, it sends $\{\langle g_k \rangle_2\}_{k \in [|\mathrm{inp}_i|]}$ to the client. The client reconstructs $\{g_k\}_{k \in [[\mathsf{inp}_i]]}$ and outputs it.

Figure 6: Our protocol for conditional disclosure of secrets.

Server's input: $\mathbf{M}_i(\mathbf{x}_i - \mathbf{r}_i) + \mathbf{s}_i, \mathbf{s}'_{i+1}$. Client's input: \mathbf{r}_{i+1} , $\mathbf{M}_i(\mathbf{r}_i) - \mathbf{s}_i$

- Compute $\mathbf{M}_i(\mathbf{x}_i) = \mathbf{M}_i(\mathbf{x}_i \mathbf{r}_i) + \mathbf{s}_i + \mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i$.
- 2. Compute $ReLU(\mathbf{M}_i(\mathbf{x}_i))$ to obtain \mathbf{x}_{i+1} .
- 3. Output $\mathbf{x}_{i+1} \mathbf{r}_{i+1} + \mathbf{s}'_{i+1}$.

Figure 7: Description of circuit C_i .

We slightly abuse the notation and use $|inp_i|$ to also denote the size of the server input to the garbled circuit.

Correctness follows from inspection; see the full version of this paper for the security proof.

Remark 5.2 (ACG for subsequent linear layers). Many network architectures contain consecutive linear layers between two ReLU activations. For simplicity of exposition, we have composed these linear layers in our protocol description. However, doing so in the actual implementation would be inefficient since our homomorphic algorithms are highly specialized for specific layer types. As a result, in practice Π_{ACG} must be modified so that on consecutive linear layers the client only receives MAC shares of the layer output on the *final* linear layer in the sequence. In the online phase, the client additionally sends MAC shares of their input on intermediate linear layers since they are not checked inside of the CDS.

Remark 5.3 (Checking client CDS inputs). f_{CDS} must check that the client's bit decomposition is correct. That is, it must check that the claimed bit decomposition (a) consists of boolean values, (b) corresponds to an integer with value less than p. For efficiency reasons, we perform only the first check in the preprocessing phase, and move the second check to our garbled circuits in the online phase.

Remark 5.4 (Fixed-point arithmetic in finite fields). Neural networks work over the real numbers, but our cryptographic protocols work over finite prime fields. To emulate real arithmetic, we rely on fixed-point arithmetic. However, to maintain precision, one needs to occasionally *truncate* intermediate values to ensure that the result does not wrap around the field.

In DELPHI, both parties perform truncation directly on their local secret shares following the technique of [Moh+17] which correctly truncates the shared value with a small, additional error. While this error does not greatly impact accuracy, it is unacceptable in the client-malicious setting as it would invalidate the MAC of the share. As a result, MUSE must perform trunction directly on the shared value using a secure MPC. We perform this truncation for free within our garbled circuits by always returning zero labels for the upper bits of the ReLU output.

5.3 An efficient protocol for computing f_{CDS}

To securely compute the function f_{CDS} , MUSE adapts the state-of-the-art arithmetic MPC framework Overdrive [Kel+18] (which achieves malicious security) to the simpler client-malicious 2PC setting. Doing so results in great efficiency improvements, as we now explain.⁸

The heaviest cryptographic costs when using Overdrive for $f_{\rm CDS}$ are due to (a) MAC key generation, (b) triples generation, and (c) authentication of secret client and server inputs. We now describe how we optimize all of these procedures in the client-malicious setting.

MAC key generation. In Overdrive, a MAC key must be secret-shared among the parties, since any party may be malicious and could use knowledge of the key to cheat. In the client-malicious setting, the server will never cheat so they can simply generate and hold the MAC key themselves.

Triples generation. In order to generate multiplication triples in Overdrive, all parties must generate ciphertexts of their shares, prove knowledge of these ciphertexts in zero-knowledge, homomorphically compute a triple from the ciphertexts, and run a distributed decryption algorithm so all parties receive a share of the result. Note that the distributed decryption allows a malicious adversary to inject an authenticated additive shift, so parties must "sacrifice" a triple in order to ensure correctness [Dam+12], which harms performance.

In the client-malicious setting, we can avoid distributed decryption, triple sacrifice, and a number of zero-knowledge proofs by taking advantage of the fact that the server knows the MAC key. In particular, we devise the following efficient protocol: the client sends the encryption of their shares directly to the server (along with a zero-knowledge proof of plaintext knowledge). The server homomorphically computes the shares of the triple, and returns it to the client. Since the server performs the computation, correctness is guaranteed and no distributed decryption or triple sacrifice is necessary. We provide benchmarks of our optimized generation in Section 6.5. See Fig. 17 for a full description of Π_{Triple} .

Input authentication. Overdrive [Kel+18] optimizes the input sharing method of [Dam+12], by assuming that the encryption scheme they employ achieves *linear-targeted malleability* (LTM) [Bit+13]. The LTM assumption for an encryption scheme informally states that only affine transformations can be computed on ciphertexts. This assumption is non-falsifiable, and, when applied to the encryption schemes used in Overdrive, has received insufficient scrutiny.

In our protocol, we avoid relying on this strong assumption by observing that the majority of secret inputs originate with the server, and because the server holds the MAC keys, it can easily authenticate its inputs *without cryptography*. In more detail, the protocol proceeds as follows.

- The server shares their inputs by producing a random authenticated share of their input using the MAC key and sends it to the client.
- The client shares their input by following the same methodology as [Dam+12]. Note that generating random authenticated shares can be implemented using our triple generation procedure from above, thus inheriting the same speedups.

We benchmark these techniques in Section 6.5. See Fig. 14 for a full description of $\Pi_{InputAuth}$.

6 Evaluation

We divide the evaluation into three sections which answer the following questions:

• Section 6.3: What are the latency and communication costs

⁸While the protocol of [Che+20] offers better performance than Overdrive, at the time of writing the source code for it was unavailable, and so we could not build upon it. Our optimizations in this section apply also to the [Che+20] protocol, so it is plausible that in the future MUSE could instead rely on it.

					MNIST		CIFAR-10	
			system	threads	time (s)	comm. (GB)	time (s)	comm. (GB)
	Linear	CG	DELPHI	1	3.93	0.03	36.21	0.05
		ACG	MUSE	1	4.74	0.04	40.78	0.07
≥00								
ssin	Non-linear	Garbling	DELPHI	2	1.81	0.18	19.31	2.95
ž		Garbling	MUSE	2	4.34	0.51	62.19	7.45
Preprocessing		OT	DELPHI	8	1.67	0.02	5.2	0.35
P		CDS Triple Gen.	MUSE	6	7.32	2.66	112.51	44.36
	~	CDS Input Auth.	MUSE	2	4.10	0.43	59.579	7.00
		CDS Evaluation	MUSE	2	2.17	0.53	31.34	8.79
Online		Online	DELPHI	8	0.48	0.01	3.74	0.16
		Online	MUSE	8	0.80	0.01	8.37	0.23

Table 3: Latency and communication cost of the individual components of MUSE and DELPHI. See Section 6.2 for more information on the network architectures and number of threads used.

of MUSE's individual components when performing inference and how do they compare to the semi-honest DELPHI?

- Section 6.4: How does MUSE compare to other inference protocols secure against malicious clients?
- Section 6.5: How do our client-malicious Overdrive subprotocols compare to standard Overdrive?

6.1 System implementation

We implemented MUSE in Rust and C++. We use the SEAL homomorphic encryption library [Sea] to implement HE, the fancy-garbling library to implement garbled circuits, and MP-SPDZ¹⁰ [Kel20] to implement zero-knowledge proofs. MUSE achieves 128 bits of computational security, and 40 bits of statistical security.

6.2 Evaluation setup

All experiments were carried out on AWS c5.9xlarge instances possessing an Intel Xeon 8000 series CPU at 3.6 GHz. The client and server instances were located in the us-west-1 (Northern California) and us-west-2 (Oregon) regions respectively with 21 ms round-trip latency. The client and server executions used 8 threads each. We evaluate MUSE on the following datasets and network architectures:

- 1. *MNIST* is a standardized dataset consisting of (28 × 28) greyscale images of the digits 0–9. The training set contains 60,000 images, while the test set has 10,000 images. Our experiments use the 2-layer CNN architecture specified in MiniONN [Liu+17a] with average pooling in place of max pooling.
- 2. CIFAR-10 is a standardized dataset consisting of (32 × 32) RGB images separated into 10 classes. The training set contains 50,000 images, while the test set has 10,000 images. Our experiments use the 7-layer CNN architecture specified in MiniONN [Liu+17a].

In our experiments, MUSE runs all of its various preprocessing components in parallel using a work-stealing thread-pool with 8 threads. For simplicity, in Table 3 we provide

microbenchmarks for each component using a static thread allocation that closely reflects the allocation used during actual execution. The current end-to-end numbers for MUSE are estimates as we have implemented all of the individual components, but are still in the process of integrating them into a full system. We carefully tested CPU time, memory usage, and bandwidth usage of each component to ensure that our estimated end-to-end numbers are accurate.

Baselines. Since there are no specialized protocols for client-malicious secure inference, we chose to use generic MPC frameworks as our baselines to compare MUSE against: maliciously-secure Overdrive [Kel+18] *and* Overdrive with our client-malicious optimizations. We used MP-SDPZ's implementation of the maliciously-secure Overdrive protocol, and estimated the total runtime and communication costs of client-malicious Overdrive using microbenchmarks from MUSE's triple generation, MUSE's input authentication, and MP-SPDZ.

Additionally, we use microbenchmarks in Table 3 to demonstrate the concrete costs of strengthening each individual component of MUSE from semi-honest to client-malicious security. As a semi-honest baseline, we chose to compare against DELPHI and not against more recent works which offer better performance [Rat+20] because (1) these newer works use different techniques for both linear and non-linear layers, which would make isolating the cost of upgrading security difficult, and (2) it is unclear how to upgrade their semi-honest protocols to achieve client-malicious security in an efficient way.

6.3 Microbenchmarks

In Table 3 we compare microbenchmarks for MUSE and DEL-PHI on the MNIST and CIFAR-10 networks using a similar number of threads to demonstrate the concrete costs of strengthening each component of DELPHI to client-malicious security.

6.3.1 Preprocessing phase

The primary difference between MUSE and DELPHI occurs in the preprocessing phase.

⁹https://github.com/GaloisInc/fancy-garbling/

 $^{^{10} {\}tt https://github.com/data61/MP-SPDZ}$

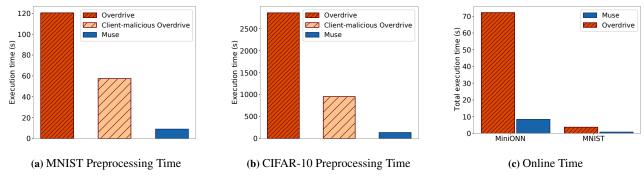


Figure 8: Comparison of execution times between MUSE, Overdrive, and client-malicious Overdrive

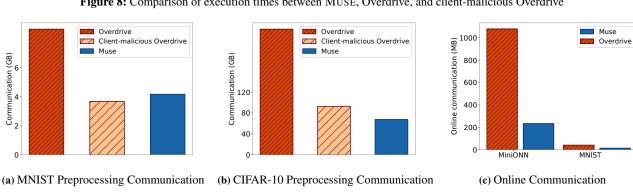


Figure 9: Comparison of communication cost between MUSE, Overdrive, and client-malicious Overdrive

Linear layers. As discussed in Section 5.1, the primary difference in how DELPHI and MUSE preprocess linear layers lies in the fact that the former uses a plain correlations generator (CG), while MUSE use an authenticated correlations generator (ACG). Because the ACG requires additional homomorphic operations and zero-knowledge proofs, we should expect MUSE to be slightly slower than DELPHI and require slightly more communication. In Table 3 we observe that this is precisely the case.

Non-linear layers. To preprocess the non-linear layers in DELPHI, the server garbles a circuit corresponding to ReLU and sends to the client. The two parties than engage in an oblivious transfer whereby the client learns the garbled labels corresponding to their input.

In MUSE, a number of modifications to this procedure must be made. First, MUSE cannot use simple oblivious transfer and must opt for the much more expensive CDS protocol to ensure the client receives the correct garbled labels. Second, as detailed in Remark 5.3, MUSE pushes some checks from the CDS to the online garbled circuits which roughly doubles the number of AND gates in the circuit.

As a result, we should expect a $2\times-3\times$ increase in latency and communication for the garbling in MUSE when compared to DELPHI, and a much higher cost for the CDS compared to oblivious transfer. Table 3 validates these hypotheses.

6.3.2 Online phase

MUSE retains the same structure for the online phase, but has a few small additions. For subsequent linear layers, MUSE

requires the client to send additional MAC shares (see Remark 5.2). For non-linear layers, MUSE requires an extra hash key to be sent, and the circuit being evaluated is roughly twice the size as the one in DELPHI.

As a result, we should expect the garbled circuit evaluation time to be the only significant difference in online runtime between MUSE and DELPHI, and for MUSE to have slightly higher communication. In Table 3, we see that the difference in online runtime is $1.7 \times -2.2 \times$ and the communication difference is approximately $1.4\times$.

In conclusion, MUSE's overhead when compared to DEL-PHI is minimal in every component except the CDS. MUSE's online phase outperforms all prior two-party semi-honest works listed in Table 1 besides DELPHI, CrypTFlow2 [Rat+20], and XONN [Ria+19].

6.4 Full system comparisons

Fig. 8 and Fig. 9 demonstrate how MUSE performs against malicious Overdrive and client-malicious Overdrive. Note that our client-malicious optimizations for Overdrive don't affect the online phase which is why we exclude client-malicious Overdrive in the online figures.

In summary, MUSE's preprocessing is $13.4 \times -21 \times$ faster and reduces communication by $2\times-3.6\times$ compared to standard Overdrive. For client-malicious Overdrive, MUSE's preprocessing phase is $6.4 \times -7 \times$ faster. For the smaller MNIST network, the communication cost of MUSE is slightly higher than that of client-malicious Overdrive (due to a constant

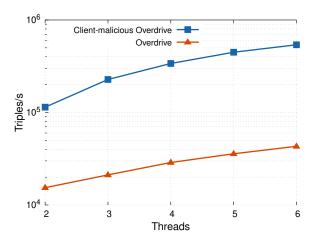


Figure 10: Triple Generation amortized on a batch of 10,000,000 triples over a 44 bit prime field with 40 bit statistical security.

	Clie	ent Inputs	Server Inputs		
Threat Model	time (s)	comm. (MB)	time (s)	comm. (MB)	
Malicious	12.11	90	12.11	90	
Client-mal.	7.406	320	0.32	20	

Table 4: Input Authentication on 1,000,000 inputs over a 44 bit prime field with 40 bit statistical security using a single thread. ■ Relies on the LTME Assumption.

overhead from the garbled circuits), but our techniques scale better and achieve a $1.4\times$ reduction for the larger CIFAR-10 network.

For the online phase, we observe a $7.8 \times -8.6 \times$ latency improvement and $3.4 \times -4.6 \times$ communication improvement when comparing MUSE to Overdrive.

6.5 Improvements to Overdrive

In this section we demonstrate the effectiveness of our optimizations to Overdrive in the client-malicious setting. In particular, we show that in client-malicious Overdrive *without the LTME assumption*:

- Triple generation is significantly more efficient.
- Client input authentication is slightly more efficient.
- Server input authentication is significantly more efficient. These improvements are of independent interest and can easily be extended to support more parties.

Triple generation. In Fig. 10 we benchmark the generation of triples on a variable number of threads. In summary, client-malicious Overdrive achieves a $8\times-12.5\times$ latency improvement and $1.7\times$ communication reduction (the latter is not shown in the graph) over standard Overdrive.

Input authentication. In Table 4 we show benchmarks for input authentication for the client and server. Our protocol for client inputs achieves a $1.6\times$ speed improvement *without the LTME assumption*, but increases communication by $3.6\times$. We observe a $37.8\times$ improvement in latency and $4.5\times$ improvement in communication for server inputs.

7 Related work

7.1 Model extraction attacks

A number of recent works extract convolutional neural networks (CNNs) [Tra+16; Mil+19; Jag+20; Rol+20; Car+20] given oracle access to a neural network. Unlike our attack, these works do not exploit properties of any secure inference protocol (and indeed do not rely on the existence of these), but require a much larger number of queries. We compare against the state-of-the-art attack [Car+20] in Section 2.3.

Mitigations. While MUSE protects against our attack, it does *not* defend against attacks that leverage only the prediction result. Mitigations fall into two camps: those that inspect prediction queries [Kes+18; Juu+19], and those that try to instead modify the network to make it resilient to extraction [Tra+16; Lee+19]. While the latter kind of defense can be applied independently of secure inference, adapting the first kind of defense to work with secure inference protocols is tricky, because it requires inspecting the client's queries, which can violate privacy guarantees.

7.2 Secure inference protocols

A number of recent works have attempted to design specialized protocols for performing secure inference. These protocols achieve efficiency by combining secure computation techniques such as homomorphic encryption [Gen09b], Yao's garbled circuits [Yao86], and homomorphic secret sharing [Boy+17] with various modifications such as approximating ReLU activations with low-degree polynomials or binarizing (quantizing to one bit of accuracy) network weights. See Table 1 for a high-level overview of these protocols.

While these works have improved on latency and communication costs by orders of magnitude, *all* of the two-party protocols in Table 1 assume a semi-honest adversary. Currently-existing maliciously-secure inference protocols generally fall into the following categories: 3PC-based protocols, generic MPC frameworks, TEE-based protocols, and GC-based protocols. In the remainder of the section we discuss each of these categories:

3PC-based protocols. Recent works have explored how the addition of a third party can greatly improve efficiency for secure machine learning applications [Moh+17; Ria+18; Moh+18; Wag+19; Wag+21; Kum+20]. Many of these protocols also allow for easy extensions to handle malicious adversaries [Moh+18; Wag+19; Wag+21]. These extensions are made possible by the fact that these works assume only *one* of the parties is corrupted. In other words, these works consider *honest majority* malicious security. On the other hand, MUSE addresses the fundamentally more difficult problem of a dishonest majority. While having three non-colluding parties is convenient from a protocol design perspective, in practice, it is difficult to setup such a third party running in a separate trust domain out of the control of the server or client.

Chameleon [Ria+18] proposed a slightly weaker threat model where a semi-honest third server assists in the preprocessing phase but is not needed for the online phase. If such a setup is feasible, MUSE could naturally take advantage of this threat model by having the semi-honest third server assist in triple generation for the CDS protocol. This augmentation improves latency and bandwidth of MUSE's preprocessing phase by roughly $3\times$.

TEE-based protocols. Generally speaking, TEE-based protocols [Tra+19; Top+18; Han+18; App19] provide better efficiency than protocols relying on purely cryptographic techniques. However, this improved efficiency comes at the cost of a weaker threat model that requires trust in hardware vendors and the implementation of the enclave. Indeed, the past few years have seen a number of powerful side-channel attacks [Bra+17; Häh+17; Göt+17; Mog+17; Sch+17; Wan+17; Van+18] against popular enclaves like Intel SGX and ARM TrustZone.

Generic frameworks. Maliciously-secure MPC frameworks exist for computing arithmetic circuits [Dam+12; Kel+18; Che+20], binary circuits [Kat+18], and mixed circuits [Rot+19; Esc+20; Moh+18]. Before MUSE, these were the only existing cryptographic mechanisms for two-party client-malicious secure inference. While [Che+20] is the most efficient of these for inference, an implementation was not available at the time of writing so we compared against [Kel+18] in Section 6.4. From the results of Section 6.4 and the experiments provided in [Che+20], we can roughly estimate that the preprocessing communication of [Che+20] is similar to MUSE, but MUSE is superior on all other accounts.

GC-based protocols. DeepSecure [Rou+18], the protocol of Ball et al. [Bal+19], and XONN [Ria+19], all use circuit garbling schemes to implement constant-round secure inference protocols. While DeepSecure supports general neural networks, the protocol of [Bal+19] operates on discretized neural networks, which have integer weights, while XONN is optimized for binarized neural networks [Cou+15], which have boolean weights. These quantized networks allow for improved performance by avoiding computing expensive fixed-point multiplication in favor of integer multiplication or binary XNOR gates.

While neural network inference is commonly performed on quantized networks [Kri18], in practice quantization is never done below 8-bits since inference accuracy begins to suffer [Ban+18]. To combat this accuracy drop, XONN increases the number of neurons in its linear layers, gaining increased accuracy at the cost of a slower evaluation time. While this technique appears to work well for the datasets XONN evaluates, additional techniques are needed to scale to more difficult datasets like Imagenet as the current best-known quantization techniques for Imagenet requires 2 bit weights and 4 bit activations [Don+19]. Consequently, it is our opinion that it is still important to focus on supporting se-

cure inference for general neural networks even though BNNs appear promising.

Any GC-based protocol can be upgraded to malicious security through a combination of cut-and-choose techniques [Zhu+16] and malicious OT-extension [Kel+15], and clientmalicious security for the evaluator by using malicious OTextension. Thus, it would follow that all of these GC-based inference protocols can be transformed into malicious and fixedsubset malicious protocols. Note that DeepSecure would provide server-malicious security since the client garbles the circuit. XONN uses a specialized protocol to evaluate the first layer of the network since the client's input is an un-quantized integer. In order for these malicious/client-malicious transformations to work, XONN would need to evaluate this layer within the more-expensive garbled circuit, instead of their optimized protocol. Furthermore, an implementation of XONN was not available at the time of writing, so we could not benchmark a client-malicious version of their protocol on our experimental setup. Finally, MUSE's online speed is already superior to the *semi-honest* versions of DeepSecure and the protocol of [Bal+19].

8 Conclusion

In this paper, we introduce a novel model-extraction attack against many semi-honest secure inference protocols which outperforms existing attacks by orders of magnitude. In response, we design and implement MUSE, an efficient two-party secure inference protocol resilient to *malicious clients*. MUSE achieves online performance close to existing semi-honest protocols, and greatly outperforms alternate solutions for client-malicious secure inference. As part of MUSE's design, we introduce a novel cryptographic protocol for *conditional disclosure of secrets* and improved procedures for generic MPC in the client-malicious setting. We hope that MUSE is a first step towards achieving practical two-party secure inference in a strong threat model.

Acknowledgements. We thank Marcel Keller for help in using MP-SPDZ, Vinod Vaikuntanathan and David Wu for answering questions about linear targeted malleable encryption, Joey Gonzalez for answering questions about Binarized Neural Networks, and the anonymous reviewers as well as our shepherd Florian Tramér for their detailed feedback. This work was supported in part by the NSF CISE Expeditions CCF-1730628, NSF Career 1943347, and gifts/awards from the Sloan Foundation, Bakar Program, Alibaba, Amazon Web Services, Ant Group, Capital One, Ericsson, Facebook, Futurewei, Google, Intel, Microsoft, Nvidia, Scotiabank, Splunk, and VMware.

References

[App19] Apple. "iOS Security". https://www.apple.com/
business/docs/site/iOS_Security_Guide.
pdf.

- [Bal+19] M. Ball, B. Carmer, T. Malkin, M. Rosulek, and N. Schimanski. "Garbled Neural Networks are Practical". ePrint Report 2019/338.
- [Ban+18] R. Banner, I. Hubara, E. Hoffer, and D. Soudry. "Scalable methods for 8-bit training of neural networks".In: NeurIPS '18.
- [Bar18] B. Barrett. "The year Alexa grew up". https://www.wired.com/story/amazon-alexa-2018-machine-learning/.
- [Bel+12] M. Bellare, V. T. Hoang, and P. Rogaway. "Foundations of garbled circuits". In: CCS '12.
- [Bit+13] N. Bitansky, A. Chiesa, Y. Ishai, R. Ostrovsky, and O. Paneth. "Succinct Non-interactive Arguments via Linear Interactive Proofs". In: TCC '13.
- [Bou+18] F. Bourse, M. Minelli, M. Minihold, and P. Paillier. "Fast Homomorphic Evaluation of Deep Discretized Neural Networks". In: CRYPTO '18.
- [Boy+17] E. Boyle, G. Couteau, N. Gilboa, Y. Ishai, and M. Orrú. "Homomorphic Secret Sharing: Optimizations and Applications". In: CCS '17.
- [Bra+17] F. Brasser, U. Müller, A. Dmitrienko, K. Kostiainen, S. Capkun, and A. Sadeghi. "Software Grand Exposure: SGX Cache Attacks Are Practical". In: WOOT '17.
- [Bru+18] A. Brutzkus, O. Elisha, and R. Gilad-Bachrach. "Low Latency Privacy Preserving Inference". ArXiV, cs.CR 1812.10659.
- [Car+20] N. Carlini, M. Jagielski, and I. Mironov. "Cryptanalytic Extraction of Neural Network Models". In: CRYPTO '20.
- [Cha+19] N. Chandran, D. Gupta, A. Rastogi, R. Sharma, and S. Tripathi. "EzPC: Programmable and Efficient Secure Two-Party Computation for Machine Learning". In: EuroS&P '19.
- [Che+20] H. Chen, M. Kim, I. P. Razenshteyn, D. Rotaru, Y. Song, and S. Wagh. "Maliciously Secure Matrix Multiplication with Applications to Private Deep Learning". In: ASIACRYPT '20.
- [Cho+18] E. Chou, J. Beal, D. Levy, S. Yeung, A. Haque, and L. Fei-Fei. "Faster CryptoNets: Leveraging Sparsity for Real-World Encrypted Inference". ArXiV, cs.CR 1811.09953.
- [Cou+15] M. Courbariaux, Y. Bengio, and J. David. "BinaryConnect: Training Deep Neural Networks with binary weights during propagations". In: NeurIPS '18.
- [Dam+12] I. Damgård, V. Pastro, N. P. Smart, and S. Zakarias. "Multiparty Computation from Somewhat Homomorphic Encryption". In: CRYPTO '12.
- [Dat+19] R. Dathathri, O. Saarikivi, H. Chen, K. Laine, K. E. Lauter, S. Maleki, M. Musuvathi, and T. Mytkowicz. "CHET: An optimizing compiler for fully-homomorphic neural-network inferencing". In: PLDI '19.
- [Don+19] Z. Dong, Z. Yao, A. Gholami, M. Mahoney, and K. Keutzer. "HAWQ: Hessian AWare Quantiza-

- tion of Neural Networks With Mixed-Precision". In: ICCV '19.
- [Esc+20] D. Escudero, S. Ghosh, M. Keller, R. Rachuri, and P. Scholl. "Improved Primitives for MPC over Mixed Arithmetic-Binary Circuits". In: CRYPTO '20.
- [Fan+12] J. Fan and F. Vercauteren. "Somewhat Practical Fully Homomorphic Encryption". ePrint Report 2012/144.
- [Gen09a] C. Gentry. "A Fully Homomorphic Encryption Scheme". PhD thesis. Stanford University, 2009.
- [Gen09b] C. Gentry. "Fully homomorphic encryption using ideal lattices". In: STOC '09.
- [Gil+16] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing. "CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy". In: ICML '16.
- [Gol+89] S. Goldwasser, S. Micali, and C. Rackoff. "The Knowledge Complexity of Interactive Proof Systems". In: SIAM J. Comput. (1989).
- [Goo17] Google. Google Infrastructure Security Design Overview. Tech. rep. 2017.
- [Göt+17] J. Götzfried, M. Eckert, S. Schinzel, and T. Müller. "Cache Attacks on Intel SGX". In: EUROSEC '17.
- [Häh+17] M. Hähnel, W. Cui, and M. Peinado. "High-Resolution Side Channels for Untrusted Operating Systems". In: ATC '2017.
- [Han+18] L. Hanzlik, Y. Zhang, K. Grosse, A. Salem, M. Augustin, M. Backes, and M. Fritz. "MLCapsule: Guarded Offline Deployment of Machine Learning as a Service". ArXiV, cs.CR 1808.00590.
- [Hes+17] E. Hesamifard, H. Takabi, and M. Ghasemi. "CryptoDL: Deep Neural Networks over Encrypted Data". ArXiV, cs.CR 1711.05189.
- [Jag+20] M. Jagielski, N. Carlini, D. Berthelot, A. Kurakin, and N. Papernot. "High Accuracy and High Fidelity Extraction of Neural Networks". In: USENIX Security '20.
- [Juu+19] M. Juuti, S. Szyller, S. Marchal, and N. Asokan. "PRADA: Protecting Against DNN Model Stealing Attacks". In: EuroS&P '19.
- [Juv+18] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. "GAZELLE: A Low Latency Framework for Secure Neural Network Inference". In: USENIX Security '18.
- [Kat+18] J. Katz, S. Ranellucci, M. Rosulek, and X. Wang. "Optimizing Authenticated Garbling for Faster Secure Two-Party Computation". In: CRYPTO '18.
- [Kel+15] M. Keller, E. Orsini, and P. Scholl. "Actively Secure OT Extension with Optimal Overhead". In: CRYPTO '15.
- [Kel+18] M. Keller, V. Pastro, and D. Rotaru. "Overdrive: Making SPDZ Great Again". In: EUROCRYPT '18.
- [Kel20] M. Keller. "MP-SPDZ: A Versatile Framework for Multi-Party Computation". In: CCS '20.

- [Kes+18] M. Kesarwani, B. Mukhoty, V. Arya, and S. Mehta. "Model Extraction Warning in MLaaS Paradigm". In: ACSAC '18.
- [Kri18] R. Krishnamoorthi. "Quantizing deep convolutional networks for efficient inference: A whitepaper". arXiv: 1806.08342 [cs.LG].
- [Kum+20] N. Kumar, M. Rathee, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. "CrypTFlow: Secure TensorFlow Inference". In: S&P '20.
- [Lee+19] T. Lee, B. Edwards, I. Molloy, and D. Su. "Defending Against Neural Network Model Stealing Attacks Using Deceptive Perturbations". In: SP Workshop '19.
- [Lin+15] Y. Lindel and P. Benny. "An Efficient Protocol for Secure Two-Party Computation in the Presence of Malicious Adversaries". In: J. Cryptol. (2015).
- [Liu+17a] J. Liu, M. Juuti, Y. Lu, and N. Asokan. "Oblivious Neural Network Predictions via MiniONN Transformations". In: CCS '17.
- [Liu+17b] W. Liu, Z. Wang, X. Liu, N. Zeng, Y. Liu, and F. E. Alsaadi. "A survey of deep neural network architectures and their applications". In: *Neurocomputing* (2017).
- [Lou+19] Q. Lou and L. Jiang. "SHE: A Fast and Accurate Deep Neural Network for Encrypted Data". ArXiV, cs.CR 1906.00148.
- [Mil+19] S. Milli, L. Schmidt, A. D. Dragan, and M. Hardt. "Model Reconstruction from Model Explanations". In: FAT* '19.
- [Mis+20] P. Mishra, R. Lehmkuhl, A. Srinivasan, W. Zheng, and R. A. Popa. "Delphi: A Cryptographic Inference Service for Neural Networks". In: USENIX Security '20.
- [Mog+17] A. Moghimi, G. Irazoqui, and T. Eisenbarth. "CacheZoom: How SGX Amplifies the Power of Cache Attacks". In: CHES '17.
- [Moh+17] P. Mohassel and Y. Zhang. "SecureML: A System for Scalable Privacy-Preserving Machine Learning". In: S&P '17.
- [Moh+18] P. Mohassel and P. Rindal. "ABY³: A Mixed Protocol Framework for Machine Learning". In: CCS '18.
- [Rab81] M. O. Rabin. "How To Exchange Secrets with Oblivious Transfer". Harvard University Technical Report 81 (TR-81).
- [Rat+20] D. Rathee, M. Rathee, N. Kumar, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. "CrypTFlow2: Practical 2-Party Secure Inference". In: CCS '20.
- [Reg09] O. Regev. "On lattices, learning with errors, random linear codes, and cryptography". In: *JACM* (2009).
- [Ria+18] 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.
- [Ria+19] M. S. Riazi, M. Samragh, H. Chen, K. Laine, K. Lauter, and F. Koushanfar. "XONN: XNOR-based

- Oblivious Deep Neural Network Inference". In: USENIX Security '19.
- [Rol+20] D. Rolnick and K. P. Körding. "Reverse-Engineering Deep ReLU Networks". In: ICML '20.
- [Rot+19] D. Rotaru and T. Wood. "MArBled Circuits: Mixing Arithmetic and Boolean Circuits with Active Security". In: INDOCRYPT '19.
- [Rou+18] B. D. Rouhani, M. S. Riazi, and F. Koushanfar. "DeepSecure: Scalable Provably-secure Deep Learning". In: DAC '18.
- [San+18] A. Sanyal, M. Kusner, A. Gascón, and V. Kanade. "TAPAS: Tricks to Accelerate (encrypted) Prediction As a Service". In: ICML '18.
- [Sch+17] M. Schwarz, S. Weiser, D. Gruss, C. Maurice, and S. Mangard. "Malware Guard Extension: Using SGX to Conceal Cache Attacks". In: DIMVA '17.
- [Sea] "Microsoft SEAL (release 3.3)". https://github.com/Microsoft/SEAL. Microsoft Research, Redmond. WA.
- [Top+18] S. Tople, K. Grover, S. Shinde, R. Bhagwan, and R. Ramjee. "Privado: Practical and Secure DNN Inference". ArXiV, cs.CR 1810.00602.
- [Tra+16] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart. "Stealing Machine Learning Models via Prediction APIs". In: USENIX Security '16.
- [Tra+19] F. Tramer and D. Boneh. "Slalom: Fast, Verifiable and Private Execution of Neural Networks in Trusted Hardware". In: ICLR '19.
- [Van+18] J. Van Bulck, M. Minkin, O. Weisse, D. Genkin, B. Kasikci, F. Piessens, M. Silberstein, T. F. Wenisch, Y. Yarom, and R. Strackx. "Foreshadow: Extracting the Keys to the Intel SGX Kingdom with Transient Out-of-Order Execution". In: USENIX Security '18.
- [Wag+19] S. Wagh, D. Gupta, and N. Chandran. "SecureNN: 3-Party Secure Computation for Neural Network Training". In: *Proc. Priv. Enhancing Technol.* (2019).
- [Wag+21] S. Wagh, S. Tople, F. Benhamouda, E. Kushilevitz, P. Mittal, and T. Rabin. "Falcon: Honest-Majority Maliciously Secure Framework for Private Deep Learning". In: *Proc. Priv. Enhancing Technol.* (2021).
- [Wan+17] W. Wang, G. Chen, X. Pan, Y. Zhang, X. Wang, V. Bindschaedler, H. Tang, and C. A. Gunter. "Leaky Cauldron on the Dark Land: Understanding Memory Side-Channel Hazards in SGX". In: CCS '17.
- [Yao86] A. C. Yao. "How to Generate and Exchange Secrets (Extended Abstract)". In: FOCS '86.
- [Zhu+16] R. Zhu, Y. Huang, J. Katz, and a. shelat. "The Cut-and-Choose Game and Its Application to Cryptographic Protocols". In: USENIX Security '16.

Functionality \mathcal{F}_{ACG}

- $\begin{array}{ll} \text{1. On client input } \{\mathbf{r}_i\}_{i \in [\ell]}, \text{ server input } \{\mathbf{s}_i, \mathbf{M}_i, \alpha_i, \beta_i\}_{i \in [\ell]}, \\ \text{ compute } \{\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i, \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i), \beta_i \cdot \mathbf{r}_i\}_{i \in [\ell]}. \end{array}$
- 2. Secret share $\alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i)$ to $\langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_1, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_2$ and $\beta_i \cdot \mathbf{r}_i$ to $\langle \beta_i \cdot \mathbf{r}_i \rangle_1, \langle \beta_i \cdot \mathbf{r}_i \rangle_2$.
- 3. Output $(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_1, \langle \beta_i \cdot \mathbf{r}_i \rangle_1)$ to the client, and $(\langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_2, \langle \beta_i \cdot \mathbf{r}_i \rangle_2)$ to the server.

Figure 11: The ideal functionality for Authenticated Correlations Generator

Functionality \mathcal{F}_{CDS}

- 1. The client and server input $(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_1, \mathbf{r}_{i+1}, \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_1)$ and $(\alpha_i, \beta_{i+1}, \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_2, \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_2, \{\mathsf{lab}_{i,k,0}^C, \mathsf{lab}_{i,k,1}^C\}_{k \in [|\mathsf{inp}_i|]})$ respectively for some $i \in [\ell]$.
- 2. If $\alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) = \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_1 + \langle \alpha_i(\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i) \rangle_2$ and $\beta_{i+1} \cdot \mathbf{r}_{i+1} = \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_1 + \langle \beta_{i+1} \cdot \mathbf{r}_{i+1} \rangle_2$, output the labels $\{\mathsf{lab}_{i,k,0}^C, \mathsf{lab}_{i,k,1}^C\}_{k \in [|\mathsf{inp}_i|]}$ corresponding to $\mathbf{M}_i(\mathbf{r}_i) \mathbf{s}_i$ and \mathbf{r}_{i+1} to the client. Otherwise, abort.

Figure 12: The ideal functionality for Conditional Disclosure of Secrets

$\mathcal{F}_{InputAuth}$

• The client's input is \mathbf{m}_c and the server's input is \mathbf{m}_s and a MAC key δ . The client receives $[\![\mathbf{m}_c]\!]_1, [\![\mathbf{m}_s]\!]_1$ and the server receives $[\![\mathbf{m}_c]\!]_2, [\![\mathbf{m}_s]\!]_2$.

Figure 13: Description of $\mathcal{F}_{InputAuth}$.

Protocol $\Pi_{InputAuth}$

- 1. Both parties invoke $\mathcal{F}_{\rm Rand}$ to receive $|\mathbf{m}_c|$ random shares $|\!|\mathbf{r}|\!|$
- 2. **r** is privately opened to the client.
- 3. The client broadcasts $\varepsilon = \mathbf{m}_c \mathbf{r}$.
- 4. The server's share is $[\![\mathbf{m}_c]\!]_2 = (\varepsilon, [\![\mathbf{r}]\!]_2)$ and the client's share is $[\![\mathbf{m}_c]\!]_1 = (\varepsilon, [\![\mathbf{r}]\!]_1)$.
- 5. The server chooses two masking vectors \mathbf{u} , \mathbf{v} and sends the client $[\![\mathbf{m}_s]\!]_1 = (\mathbf{m}_s \mathbf{u}, \mathbf{\delta} \cdot \mathbf{m}_s \mathbf{v})$. The server sets its share $[\![\mathbf{m}_s]\!]_2 = (\mathbf{u}, \mathbf{v})$

Figure 14: The protocol for Input Authentication.

$\mathcal{F}_{\text{Triple}}$

• The client's input is a_1,b_1 and the server's input is a_2,b_2 and a MAC key δ . The client receives $[\![a_1+a_2]\!]_1,[\![b_1+b_2]\!]_1,[\![(a_1+a_2)\cdot(b_1+b_2)]\!]_1$ and the server receives $[\![a_1+a_2]\!]_2,[\![b_1+b_2]\!]_2,[\![(a_1+a_2)\cdot(b_1+b_2)]\!]_2$.

Figure 15: Description of \mathcal{F}_{Triple} .

\mathcal{F}_{Rand}

• The client's input is \mathbf{r}_1 and the server's input is \mathbf{r}_2 and a MAC key δ . The client receives $[\![\mathbf{r}_1 + \mathbf{r}_2]\!]_1$ and the server receives $[\![\mathbf{r}_1 + \mathbf{r}_2]\!]_2$.

Figure 16: Description of \mathcal{F}_{Rand} .

Protocol Π_{Triple}

- i. The client and the server engage in a two-party computation protocol with security against malicious clients and semi-honest servers to generate the public key, secret key pair for HE. At the end of the protocol, the client learns the public key pk and the secret key sk whereas the server only learns pk.
- ii. The client sends $\mathsf{Enc}(\mathsf{pk}, a_1), \mathsf{Enc}(\mathsf{pk}, b_1)$ to the server along with a zero-knowledge proof of well-formedness of the two ciphertexts. The server verifies this proof before continuing.
- iii. The server homomorphically computes $\operatorname{Enc}(\operatorname{pk},a_1+a_2)$, $\operatorname{Enc}(\operatorname{pk},b_1+b_2)$ and $\operatorname{Enc}(\operatorname{pk},(a_1+a_2)\cdot(b_1+b_2))$ along with $\operatorname{Enc}(\operatorname{pk},\delta(a_1+a_2))$, $\operatorname{Enc}(\operatorname{pk},\delta(b_1+b_2))$ and $\operatorname{Enc}(\operatorname{pk},\delta\cdot(a_1+a_2)\cdot(b_1+b_2))$.
- iv. The server chooses six random masking elements $u_1, v_1, t_1, u_2, v_2, t_2$ and computes $\operatorname{Enc}(\operatorname{pk}, a_1 + a_2 u_1)$, $\operatorname{Enc}(\operatorname{pk}, b_1 + b_2 v_1)$ and $\operatorname{Enc}(\operatorname{pk}, (a_1 + a_2) \cdot (b_1 + b_2) t_1)$ along with $\operatorname{Enc}(\operatorname{pk}, \delta(a_1 + a_2) u_2)$, $\operatorname{Enc}(\operatorname{pk}, \delta(b_1 + b_2) v_2)$ and $\operatorname{Enc}(\operatorname{pk}, \delta \cdot (a_1 + a_2) \cdot (b_1 + b_2) t_2)$. It sends these six ciphertexts to the client.
- v. The client decrypts the above ciphertexts and obtains $[\![a_1+a_2]\!]_1=(a_1+a_2-u_1,\delta(a_1+a_2)-u_2),[\![b_1+b_2]\!]_1=(b_1+b_2-v_1,\delta(b_1+b_2)-v_2),[\![(a_1+a_2)\cdot(b_1+b_2)]\!]_1=(a_1+a_2)\cdot(b_1+b_2)-t_1,\delta\cdot(a_1+a_2)\cdot(b_1+b_2)-t_2).$ The server outputs $[\![a_1+a_2]\!]_2=(u_1,u_2),[\![b_1+b_2]\!]_2=(v_1,v_2),[\![(a_1+a_2)\cdot(b_1+b_2)]\!]_2=(t_1,t_2).$

Figure 17: The protocol for Triple Generation.

Pseudocode for our attacks from Section 2

RecoverNetwork:

- 1. First, recover the last layer:
 - (a) Denote by \tilde{M}_{ℓ} the recovered matrix for the last layer.
 - (b) For each $i \in [t]$:
 - i. Set the initial input to the network to be zero, i.e. $\mathbf{x}_1 := 0$.
 - ii. Follow the inference protocol to partially evaluate the network up to the $\ell-1$ -th layer:
 - $\mathbf{x}_{\ell-1} := \text{ReLU}(M_{\ell-1}(\dots \text{ReLU}(M_1\mathbf{x}_1))) = \mathbf{0}.$
 - iii. Malleate the client's share of $\mathbf{x}_{\ell-1}$: $\langle \mathbf{x}'_{\ell-1} \rangle_C :=$ $\langle \mathbf{x}_{\ell-1} \rangle_C + \mathbf{e}_i$.
 - iv. Complete the protocol with the server to obtain $\mathbf{x}_{\ell} := M_{\ell} \mathbf{x}'_{\ell-1} = M_{\ell} \mathbf{e}_{j}.$
 - v. Set the *j*-th column of \tilde{M}_{ℓ} to be \mathbf{x}_{ℓ} .
 - (c) Output \tilde{M}_{ℓ} .
- 2. Then, recover all previous layers:
 - (d) For each $i \in [\ell 1, ..., 1]$:
 - i. If the *i*-th layer is a fully-connected layer, set $M_i := \mathsf{RecoverFCLayer}(M_{i+1}, \dots, M_{\ell}).$
 - ii. If the *i*-th layer is a convolutional layer, set M_i := RecoverConvLayer($M_{i+1}, \ldots, M_{\ell}$).
 - (e) Output $(M_1, ..., M_{\ell-1})$.

RecoverConvLayer $(M_{i+1}, \ldots, M_{\ell})$:

- 1. Let the dimensions of the convolutional kernel K_i be $k_i \times k_i$.
- 2. Sample a random matrix *R* having the same dimension as \mathbf{x}_{i-1} .
- 3. Apply the im2col transformation to R to obtain R'.
- 4. Let S be the indices of the pivot columns of R' when it is in row-reduced echelon form. If $|S| < k_i \times k_i$, resample R and retry. Then S specifies the indices of the independent columns of R'.
- 5. Follow the inference protocol to evaluate the network up to the i-1-th layer to obtain (a share of) the intermediate state $\mathbf{x}_{i-1} := \text{ReLU}(M_{i-1}(...\text{ReLU}(M_1\mathbf{x}_1))) = \mathbf{0}.$
- 6. Malleate the client's share of \mathbf{x}_{i-1} : $\langle \mathbf{x}'_{i-1} \rangle_C := \langle \mathbf{x}_{i-1} \rangle_C + R$.
- 7. Interact with the server to evaluate the i-th linear layer to obtain a share of $\mathbf{y}_i := M_i \mathbf{x}'_{i-1}$.
- 8. Obtain the input for the next linear layer: $\langle \mathbf{x}_i \rangle_C :=$ MaskAndLinearizeReLU($\langle \mathbf{y}_i \rangle_C, S$).
 - The vector \mathbf{x}_i is now all-zero, except at locations in S, where it equals the corresponding elements of y_i .
- 9. Interact with the server to complete the evaluation of the rest of the network, invoking LinearizeReLU to force intervening ReLUs to behave linearly.
- 10. Set $K := [X_1, ... X_{k_i \times k_i}]$, where each X_i is a formal variable.
- 11. Set X to be the all-zero matrix of dimension equal to R', except at locations in S, where it equals R'.
- 12. Compute $\mathbf{X}_{\ell} := M_{\ell} \cdot M_{\ell-1} \cdots M_{i+1} \cdot (KX)$.
- 13. Solve the linear system $\mathbf{X}_{\ell} = \mathbf{x}_{\ell}$ to learn the values of the formal variables X_i , and hence the kernel K_i .

MaskAndLinearizeReLU($\langle \mathbf{y} \rangle_C$, S):

- 1. Malleate the client's local share of y to obtain a share of the malleated y' as follows:
 - (a) For all $l \in S$, set $\langle \mathbf{y}' \rangle_C[l] := \langle \mathbf{y} \rangle_C[l] + c$.
- (b) For all $l \notin S$, set $\langle \mathbf{y}' \rangle_C[l] := \langle \mathbf{y} \rangle_C[l] c$. 2. Obtain $\langle \mathbf{x} \rangle_C := \text{LinearizeReLU}(\mathbf{y}')$..
- 3. Invert Step 1a by malleating $\langle \mathbf{x} \rangle_C$: set $\langle \mathbf{x}' \rangle_C[S] := \langle \mathbf{x} \rangle_C[S] c$
- 4. Output $\langle \mathbf{x}' \rangle_C$.

LinearizeReLU($\langle \mathbf{y} \rangle_C$):

- 1. Malleate the client's local share of y to obtain a share of the malleated \mathbf{y}' : set $\langle \mathbf{y}' \rangle_C := \langle \mathbf{y} \rangle_C + c$.
- 2. Interact with the server to obtain $\langle \mathbf{x} \rangle_C$, which is the client's share of $\mathbf{x} := \mathsf{ReLU}(\mathbf{y}')$.
- 3. Invert Step 1 by malleating $\langle \mathbf{x} \rangle_C$: set $\langle \mathbf{x}' \rangle_C := \langle \mathbf{x} \rangle_C c$
- 4. Output $\langle \mathbf{x}' \rangle_C$.

RecoverFCLayer($M_{i+1}, \ldots, M_{\ell}$):

- 1. Let the dimension of the *i*-th linear layer be $s_i \times t_i$.
- 2. Set M'_i to be a $s_i \times t_i$ matrix consisting of formal variables.
- 3. Let $s_i' := \lfloor s_i/m \rfloor$.
- 4. For each $j \in [t_i]$, and for each $k \in [s'_i]$:
- (a) Set the initial input to the network to be zero, i.e. $x_1 := 0$.
- (b) Follow the inference protocol to evaluate the network up to the i-1-th layer to obtain (a share of) the intermediate state $\mathbf{x}_{i-1} := \mathsf{ReLU}(M_{i-1}(\dots \mathsf{ReLU}(M_1\mathbf{x}_1))) = \mathbf{0}.$
- (c) Construct a query $\mathbf{q}_i := \mathbf{e}_i$.
- (d) Malleate the client's share of \mathbf{x}_{i-1} : $\langle \mathbf{x}'_{i-1} \rangle_C := \langle \mathbf{x}_{i-1} \rangle_C + \mathbf{q}_j$.
- (e) Interact with the server to evaluate the i-th linear layer to obtain a share of $\mathbf{y}_i := M_i \mathbf{x}'_{i-1}$.
- (f) Set $k' := k \cdot m$, and $S := \{k', \dots, k' + m 1\}$.
- (g) Obtain the input for the next linear layer: $\langle \mathbf{x}_i \rangle_C := \mathsf{MaskAndLinearizeReLU}(\mathbf{y}_i, S).$ The vector \mathbf{x}_i is now all-zero, except at locations in S, where it equals the corresponding elements of y_i .
- (h) Interact with the server to complete the evaluation of the rest of the network, invoking LinearizeReLU to force intervening ReLUs to behave linearly.
- (i) Compute \mathbf{X}_i as follows. First, compute $M'_i \cdot \mathbf{q}_i$, and then zero out all locations that are not in S.
- (j) Construct the *k*-th linear system $\mathbf{x}_{\ell} = M_{\ell} \cdot M_{\ell-1} \cdots M_{i+1} \cdot \mathbf{X}_{i}$.
- 5. Solve all the linear systems to recover the matrix M_i .

^aThe pivot columns are linearly independent by definition, and row operations do not change linear dependence of columns.