

Stealing Neural Network Models through the Scan Chain: A New Threat for ML Hardware

Seetal Potluri
 Department of ECE
 North Carolina State University
 Raleigh, U.S.
 spotlur2@ncsu.edu

Aydin Aysu
 Department of ECE
 North Carolina State University
 Raleigh, U.S.
 aaysu@ncsu.edu

Abstract—Stealing trained machine learning (ML) models is a new and growing concern due to the model’s development cost. Existing work on ML model extraction either applies a mathematical attack or exploits hardware vulnerabilities such as side-channel leakage. This paper shows a new style of attack, for the first time, on ML models running on embedded devices by abusing the scan-chain infrastructure. We illustrate that having course-grained scan-chain access to non-linear layer outputs is sufficient to steal ML models. To that end, we propose a novel *small-signal analysis* inspired attack that applies small perturbations into the input signals, identifies the quiescent operating points and, selectively activates certain neurons. We then couple this with a *Linear Constraint Satisfaction* based approach to efficiently extract model parameters such as weights and biases. We conduct our attack on neural network inference topologies defined in earlier works, and we automate our attack. The results show that our attack outperforms mathematical model extraction proposed in CRYPTO 2020, USENIX 2020, and ICML 2020 by an increase in accuracy of $2^{20.7} \times$, $2^{50.7} \times$, and $2^{33.9} \times$, respectively, and a reduction in queries by $2^{6.5} \times$, $2^{4.6} \times$, and $2^{14.2} \times$, respectively.

I. INTRODUCTION

Stealing trained machine learning (ML) models is a significant challenge the industry faces given the exponential increase in model development costs [1]. Indeed, prior work has considered *model stealing/extraction* attacks in the context of ML-based cloud services [2]–[4], where the motivation is to exploit the commercial value of the ML model IP. Orthogonally, others have proposed using the extracted model to craft superior adversarial samples [5]–[7]. These works assume an ML-as-a-Service application with the prediction application programmer interface and publicly accessible query interfaces.

There is, however, a trend for edge intelligence where the trained machine learning / AI algorithms execute on embedded devices, instead of on the cloud, to improve energy cost or response time. Unfortunately, such embedded applications cause new threat vectors on model extraction like *physical* side-channel attacks [8]–[13] and fault-injection attacks [14].

In this paper, we expose a new threat vector for model extraction on edge devices: *scan-chain attacks*. These attacks exploit the scan-chain infrastructure that allows in-field testing of deployed devices [15]–[17]. Although such attacks are well-known and well-recognized for cryptographic implementations [18]–[22], they have never been explored on ML hardware before. We show that this extension is non-trivial—scan-chain access can be course-grained and neural networks have unique functions and floating-point arithmetic that do not occur in cryptography. Therefore, a successful and efficient attack requires innovative

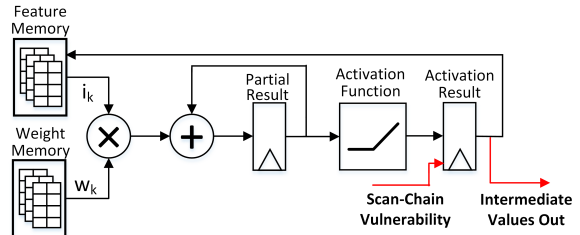


Fig. 1: Scan-chains enable accessing internal states and leaking intermediate results, which can be used to efficiently extract parameters.

techniques to steal weight and bias parameters of a neural network with a high-precision and with a low number of queries.

The contributions of this paper are as follows.

- We propose the first model stealing attack through scan-chain access. This attack exploits the *debug port* to observe the hardware internal states for model extraction.
- We show that the attack does not require accessing all internal registers. Instead, observing non-linear (i.e., activation layer or max-pooling layer) output is sufficient for the attack with two proposed novel techniques:
 - 1) Applying input perturbations and performing *small-signal analysis* to systematically expose the model parameters of each neuron.
 - 2) Formulating selective neuron activation as a linear constraint satisfaction (LCS) problem and using *coupled constraints* for *{negative weight, negative bias}* cases while using *iterative constraint relaxation* for *hard cases* of the LCS problem.
- We show the attack’s feasibility on models used in earlier works and defined for popular edge-device hardware IPs such as Eyeriss [23]. We automate the proposed attack.
- The results show that on average, our proposed attack is $2^{20.7} \times$, $2^{50.7} \times$, and $2^{33.9} \times$ more accurate and requires $2^{6.5} \times$, $2^{4.6} \times$, and $2^{14.2} \times$ fewer queries compared to the attacks proposed at CRYPTO 2020 [24], USENIX 2020 [25], and ICML 2020 [26] respectively.

II. THREAT MODEL

We follow the standard threat model laid out in the seminal work of scan-chain attacks on cryptographic circuits [18] and detailed in subsequent works [19]–[22]. We assume that the adversary can find the positions of all scan elements in the scan-chain, either by loading pairs of known plaintexts with a one-bit difference in the normal mode and then scanning out the internal state in the test mode [18] or through semiconductor reverse engineering [27].

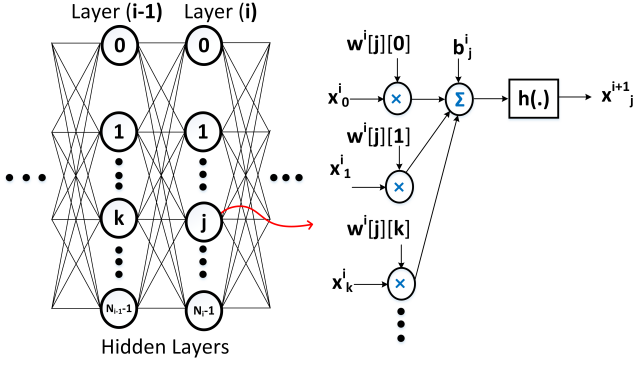


Fig. 2: The connections between two hidden layers of an FCNN

The adversary targets the inference and its goal is to extract trained model parameters such as weights and biases of a deep neural network. We consider an adversary having physical access to an edge device that performs the inference. The adversary executes the classification for a certain number of functional cycles, switches to debug mode, and uses the hardware debug port of the edge device to leak internal states as shown in Figure 1. The adversary then performs a statistical analysis on these internal states to steal the NN model parameters.

There are dedicated registers for storing the configuration bits in neural network hardware accelerators [23]. These configuration registers are tamper-proof, thus making it infeasible for the adversary to ascertain the contents of the partial results. By contrast, the contents of the ReLU activation-unit output register can be ascertained because, after a fixed number of clock cycles, it will contain the output of a certain neuron in a certain layer, before sending it to the DRAM [23]. Therefore, our attack only assumes access to the activation results shown in Figure 1.

Recent work categorizes model extraction attacks into high fidelity vs. high accuracy [25]: our adversary performs a *high fidelity* attack that aims for general agreement between the extracted and victim models on any input. Following earlier works on model extraction [4], [11], [25], we assume the hyper-parameters of the model like the types/number of neural network layers are either public or obtained by another attack [6], [12].

III. BACKGROUND

A neural network (NN) can be described as a series of functional transformations, where the architecture and parameters completely describe the NN model. A training set is used to train/tune the parameters of an adaptive NN model. Once the training is complete, the network is used to perform classification or regression. This work focuses on NN models for classification. A fully connected neural network (FCNN) is a feed-forward neural network, consisting of fully connected layers which connect every neuron in one layer to every neuron in another layer. Figure 2 shows two hidden layers of an FCNN. The outputs of neurons in $(i+1)^{th}$ ($0 \leq i < d$) fully connected layer can be described as:

$$\mathbf{x}^{i+1} = h(\mathbf{W}^i \cdot \mathbf{x}^i + \mathbf{b}^i) \quad (1)$$

where $h(\cdot)$ is the non-linear activation function, and d is the number of hidden layers. Here, $\mathbf{x}^i = [x_0^i \ x_1^i \ \dots \ x_k^i \ \dots \ x_{(N_i-1)}^i]^T$ (column vector) is the input vector to layer- (i) , where each variable x_k^i corresponds to the $(k+1)^{th}$ input feature to the layer and

$$D = \begin{cases} N_I & i=0 \\ N_{(i-1)} & i \neq 0 \end{cases}$$

where N_I is the input dimension and $N_{(i-1)}$ is the number of neurons in layer- $(i-1)$ (prior layer). Similarly, $\mathbf{x}^{i+1} = [x_0^{i+1} \ x_1^{i+1} \ \dots \ x_j^{i+1} \ \dots \ x_{(N_i-1)}^{i+1}]^T$ (column vector) is the output vector of layer- (i) , where N_i is the number of neurons in layer- (i) (current layer) and each variable x_j^{i+1} corresponds to the output of neuron- (j) in layer- (i) . The secret model parameters are the weight matrix coefficients, where $\mathbf{W}^i = \{w^i[j][k]\}$ is the weight matrix of layer- (i) , whose $(j+1)^{th}$ row corresponds to the weight vector of neuron- (j) in layer- (i) and bias coefficients, where $\mathbf{b}^i = [b_0^i \ b_1^i \ \dots \ b_j^i \ \dots \ b_{(N_i-1)}^i]^T$ (column vector) is the bias vector of layer- (i) , where b_j^i is the bias value of neuron- (j) in layer- (i) , as shown in Figure 2. These weights and biases are obtained during training and is then constant during inference. The weights, biases, and activation values are represented as double-precision floating point numbers in this work.

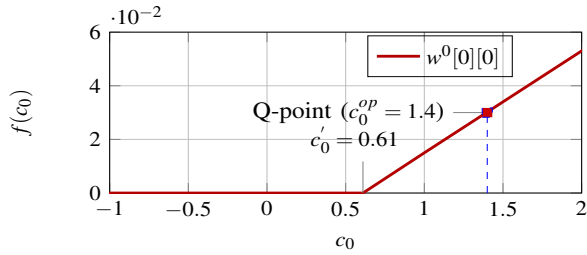
The rectified linear unit (ReLU) [28] is the non-linear activation function $h(\cdot)$ we used. This is due to ReLU's widespread success and to be consistent with the earlier works on cryptanalytic model extraction [24], [25]. The ReLU activation function is given below where x is the output of a neuron in a fully connected layer and $h(x)$ is the activation result:

$$h(x) = \max(0, x) = \begin{cases} 0 & x < 0 \\ x & x \geq 0 \end{cases}$$

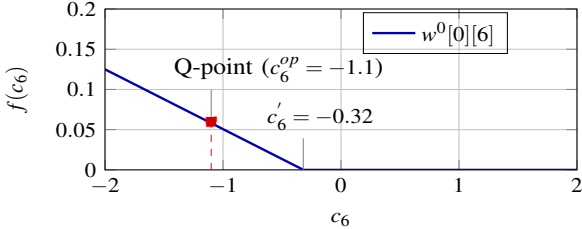
The output layer uses *softmax* operation for generating the confidence scores and subsequently selecting the output class with the maximum confidence score. During inference, the output layer performs an order-preserving transformation on the output signals of the penultimate layer [24], [25] and uses a $\max(\cdot)$ operation over these transformations, and selects the corresponding index as the output of the classifier.

IV. THE PROPOSED ATTACK

To steal NN model parameters, we apply an attack inspired by *small-signal analysis* (SSA) used to extract analog system model parameters like gain and transconductance [29]. We use SSA to identify the linear region of operations for neurons where the ReLU returns a non-zero value. We then use this to apply small perturbations into neural network input signals to generate a system of linear constraints (SLC) with known inputs and unknown parameters and solve the system using a *Linear Constraint Solver*. We address the challenges of exposing multiple weights and the parameters in subsequent layers through an iterative procedure. The iterative procedure (a) uses the extracted parameters of the previous hidden layers to generate the SLC for the current hidden layer, (b) solve the SLC to identify the neuron quiescent operating point (*Q-point*), (c) apply small perturbations into neural network input signals around the *Q-point*, let them propagate until the current hidden layer, and (d) use the black-box responses of the current hidden layer to extract the model parameters through SSA. We finally automate this iterative procedure to extend the proposed attack to FCNNs of arbitrary topologies and depths.



(a) Neuron output signal when attacking its 1st filter weight.



(b) Neuron output signal when attacking its 7th filter weight.

Fig. 3: Varying neuron signal behaviors, depending on the **weight-under-attack**. Please note the Y-axis range.

A. Identifying neuron quiescent operating point (*Q-point*)

The goal of *Q-point* identification is to search for an input that activates desired neurons and deactivates the remaining neurons in a stable way. In this context, stability refers to the activated neurons remaining active and the inactive neurons remaining inactive even in the presence of small-signal perturbations at the input. ReLU makes inactive neuron to always produce a “0” output, which does not provide useful information about the NN model. By contrast, an activated neuron’s output is a linear function of its model parameters, which leaks information and can be extracted. Such carefully selected inputs, popularly known quiescent operating points (*Q-points*) in the analog signal processing literature, are coupled with *SSA* to extract circuit model parameters. Our model extraction is inspired by this strategy. This section describes *Q-point* identification/search.

The weighted summation operation of a neuron can be considered as a filtering operation, and the array of weights can be considered as a neural filter [23]. Suppose we are interested in the output signal of the first neuron in the first hidden layer of an FCNN with the filter defined as:

$$\mathbf{W}^0[0] = \begin{bmatrix} w^0[0][0] & w^0[0][1] & \dots & w^0[0][N_I - 1] \end{bmatrix}$$

If we are interested to attack $w^0[0][0]$ (the superscript “0” denotes the first hidden layer, the first index “0” indicates the first neuron in the layer, and the second index “0” indicates the first weight in the filter), we can make all the entries in the input vector \mathbf{I} zero except the first entry containing a user-defined variable c_0 , represented as:

$$\mathbf{I} = \begin{bmatrix} c_0 & 0 & \dots & 0 \end{bmatrix}$$

We refer to the classifier as an *oracle*, to which we make queries and record the corresponding internal states. If we make a query to the oracle by applying this input and observe output of first neuron in the first layer ($O^0[0]$) through the debug port, we notice that the signal is purely a function of c_0 :

$$f(c_0) = O^0[0] = \text{ReLU}(c_0 \cdot w^0[0][0] + b^0[0]) \quad (2)$$

Figure 3(a) shows the response of neuron-(0) in layer-(0) for varying input (c_0). The example reflects the case where an

FCNN trained on MNIST and mapped to an accelerator (e.g., Eyeriss [23]), where our attack makes oracle queries with different values of c_0 and observes the debug port output. We make the following observations from this signal characteristic:

- The neuron output signal is “0” for $c_0 < 0.61$ and behaves as a linear function of the weight $w^0[0][0]$ with any further increase in c_0 . This threshold is denoted as c_0' .
- Positive perturbations to c_0 around c_0' activate the neuron, while negative perturbations do not. Thus, c_0 is sensitive around c_0' and the operating point of c_0 should be set to a value reasonably higher than c_0' so that small perturbations either way (positive or negative) do not disturb the linear behavior of the neuron.

If we denote the operating point of c_0 as c_0^{op} , then how greater it has to be as compared to c_0' depends on the magnitude of the small-signal. If we denote small-signal on c_0 as δ_0 and the amplitude of δ_0 as Δ , then we can safely define the operating point as $c_0^{op} = c_0' + \beta \cdot \Delta$, where $\beta > 1$, is a tuning parameter used to adjust the *Q-point* to the stable region of operation (i.e., both positive and negative perturbations with amplitude Δ on c_0 around the *Q-point* are acceptable). By this definition of operating point, we know for sure that $f(c_0) = f(c_0^{op} \pm \delta_0) > 0$, hence forcing the neuron to operate in the linear region of operation. Although we have so far explained the procedure to extract $w^0[0][0]$, the methodology is clearly generic. Figures 3(a) and 3(b) show the stable neuron operating points in the linear region, when attacking weights $w^0[0][0]$ (1st weight in the filter of neuron-(0) in layer-(0)) and $w^0[0][6]$ (7th weight in the filter of neuron-(0) in layer-(0)) respectively.

B. Small-Signal Analysis (*SSA*)

The connection between two neurons have two unknown variables: *weight vector* and *bias*, where *bias* is fixed for all the connections to a neuron in the current hidden layer. Therefore, obtaining two linear equations is sufficient to solve and extract an *unknown weight* for a *given neuron*. After computing a stable neuron operating point (*Q-point*) in the linear region of operation, we achieve this with an *SSA* that consist of primarily two phases:

- 1) Applying input \mathbf{I} to the neuron (since it is the first hidden layer, neuron \mathbf{I} is directly applied to the neuron input; for the subsequent hidden layers, a transformed version of \mathbf{I} gets applied) and computing its output.
- 2) Applying a small perturbation $\delta < \Delta$ to \mathbf{I} to produce \mathbf{I}^p , applying \mathbf{I}^p to the neuron and computing its output.

The results computed in these two stages will be useful for the extraction of model weights and biases. In our example, we obtain $f(c_0^{op})$ and $f(c_0^{op} + \delta_0)$ in these two phases. The difference between these two signals gives:

$$\begin{aligned} f(c_0^{op} + \delta_0) - f(c_0^{op}) &= \text{ReLU}((c_0^{op} + \delta_0) \cdot w^0[0][0] + b^0[0]) \\ &\quad - \text{ReLU}(c_0^{op} \cdot w^0[0][0] + b^0[0]) \\ &= ((c_0^{op} + \delta_0) \cdot w^0[0][0] + b^0[0]) \\ &\quad - (c_0^{op} \cdot w^0[0][0] + b^0[0]) = \delta_0 \cdot w^0[0][0] \quad (3) \end{aligned}$$

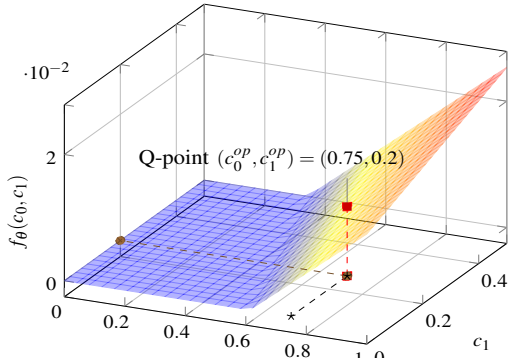


Fig. 4: $f_{\theta}(c_0, c_1)$ is a surface, and the linear region is a *hyperplane* in the 3-dimensional space.

Hence $f(c_0^{op} + \delta_0) - f(c_0^{op}) = \delta_0 \cdot w^0[0][0]$ (this has become plausible because of the linear operating point). Thus, the filter weight can be extracted as:

$$w^0[0][0] = \frac{f(c_0^{op} + \delta_0) - f(c_0^{op})}{\delta_0} \quad (4)$$

If we substitute this value for $w^0[0][0]$ in Equation 2, we obtain the bias using:

$$b^0[0] = f(c_0^{op}) - c_0^{op} \cdot w^0[0][0] = \left(1 + \frac{c_0^{op}}{\delta_0}\right) \cdot f(c_0^{op}) - \left(\frac{c_0^{op}}{\delta_0}\right) \cdot f(c_0^{op} + \delta_0) \quad (5)$$

The above calculations return $b^0[0] = -2.36e^{-2}$ and $w^0[0][0] = 3.85e^{-2}$, which are indeed the correct values for the target FCNN. Thus, we have successfully and accurately obtained both the filter weight and neuron bias. Note that this is possible given scan-chain access to the output of the non-linear layer of activation results. In our example, we have chosen the non-zero entry in location-(0) of the input signal \mathbf{I} , hence we are able to extract the 1st filter weight in neuron-(0) $w^0[0][0]$ successfully. In fact, we can choose any location-(k) in neuron-(j) and successfully extract corresponding filter weight $w^0[j][k]$, $\forall j \in [0, N_0)$, $\forall k \in [0, N_f)$. We will show later in Section IV-D, that this can indeed be generalized to other layers as well.

C. Exposing multiple filter weights using a common Q-point

So far, for simplicity, we have seen the analysis corresponding to exposing one filter weight at a time. Although this method works well for $\{\text{positive weight, positive bias}\}$ (case-1), $\{\text{positive weight, negative bias}\}$ (case-2) and $\{\text{negative weight, positive bias}\}$ (case-3) cases, this method does *not* extend to $\{\text{negative weight, negative bias}\}$ (case-4) case in the subsequent hidden layers¹. This is because all subsequent layers receive only non-negative inputs due to the ReLU activation. To address this limitation, it is important to find a *common Q-point* for multiple weights inside a neuron, and then expose the individual weights. Consider the simple example of applying the following input with only two non-zero entries at locations (0) and (1):

$$\mathbf{I} = \begin{bmatrix} c_0 & c_1 & 0 & \dots & 0 \end{bmatrix}$$

In this scenario, the output signal (of neuron-(0) in layer-(0)) is a function of both c_0 and c_1 :

$$f(c_0, c_1) = O^0[0] = \text{ReLU}(c_0 \cdot w^0[0][0] + c_1 \cdot w^0[0][1] + b^0[0]) \quad (6)$$

Since $f(c_0, c_1)$ is a rectified linear function of the 2 user-defined variables c_0 and c_1 , the neuron output function is

¹ It was possible for the case shown in Figure 3(b), because negative inputs are possible for the first hidden layer.

a surface and the linear region is a *hyperplane*, both in 3-dimensional space as shown on Figure 4. Going a step forward, if we allow as many non-zero entries as the number of filter weights at corresponding locations, then $\mathbf{I} = \mathbf{C}$ where

$$\mathbf{C} = \begin{bmatrix} c_0 & c_1 & \dots & c_{(N_f-1)} \end{bmatrix}$$

then the corresponding neuron output signal is given by:

$$f(\mathbf{C}) = O^0[0] = \text{ReLU}(\mathbf{C} \odot \mathbf{W}^0[0] + b^0[0]) \quad (7)$$

where \odot is the element-by-element (scalar) multiplication between the vectors. Here also, $f(\mathbf{C})$ is a rectified linear function of the elements in \mathbf{C} . If n be the dimension of the vector space for \mathbf{C} , then $f(\mathbf{C})$ is a surface and the linear region is a *hyperplane*, both in $(n+1)$ -dimensional space.

Let \mathbf{C}^{op} be the linear operating point (*Q-point*), and let $\mathbf{C}^p(k)$ be the perturbed vector with perturbation only in location-(k). Suppose we are interested in attacking weight $w^0[0][0]$, we can then define $\mathbf{C}^p(0)$ as:

$$\mathbf{C}^p(0) = \begin{bmatrix} c_0 + \delta_0 & c_1 & \dots & c_{(N_f-1)} \end{bmatrix}$$

In that case, the weight $w^0[0][0]$ can be extracted as follows:

$$\begin{aligned} f(\mathbf{C}^p(0)) - f(\mathbf{C}^{op}) &= \text{ReLU}(\mathbf{C}^p(0) \odot \mathbf{W}^0[0] + b^0[0]) \\ &\quad - \text{ReLU}(\mathbf{C}^{op} \odot \mathbf{W}^0[0] + b^0[0]) \\ &= (\mathbf{C}^p(0) \odot \mathbf{W}^0[0] + b^0[0]) \\ &\quad - (\mathbf{C}^{op} \odot \mathbf{W}^0[0] + b^0[0]) = w^0[0][0] \cdot \delta_0 \end{aligned} \quad (8)$$

In Equation 8, the ReLU function disappears because the neuron is operating around *Q-point*, and is thus in linear region of operation. This means the filter weight can be calculated as:

$$w^0[0][0] = \frac{f(\mathbf{C}^p(0)) - f(\mathbf{C}^{op})}{\delta_0} \quad (9)$$

All weights within the filter $\mathbf{W}^0[0]$ can be extracted similarly:

$$w^0[0][k] = \frac{f(\mathbf{C}^p(k)) - f(\mathbf{C}^{op})}{\delta_k}, 0 \leq k < D \quad (10)$$

Subsequently, the neuron bias for the neuron-(0) in layer-(0) can be extracted using:

$$b^0[0] = f(\mathbf{C}^{op}) - \mathbf{C}^{op} \odot \mathbf{W}^0[0] \quad (11)$$

This can further be generalized to all neurons in layer-(0).

D. Exposing model parameters beyond the first hidden layer

In the first layer, the entries of the input and the filter have a one-to-one mapping because of which exposing exactly one weight through SSA was possible. By contrast, the small-signal of any given index in the input \mathbf{I} gets broadcasted to all inputs to the second layer. For example, if we apply a small-signal δ_0 to c_0 , we know that $O^0[0]$ will be a function of δ_0 . Similarly, $\{\mathbf{O}^0[j]\}$, $0 \leq j < N_0$ are all functions of δ_0 , thereby each weight of a given neuron in the second layer gets multiplied by a function of δ_0 . This not only restricts direct exposing of weights like the first layer but also results in an *under-determined system* of linear equations, hence difficult to solve.

1) *Under-determined system of equations:* If we consider again the single perturbation scenario:

$$\mathbf{C}^p(0) = \begin{bmatrix} c_0 + \delta_0 & c_1 & \dots & c_{(N_f-1)} \end{bmatrix}$$

The SSA analysis on $(j+1)^{th}$ input to any second layer neuron is given as $\nabla \mathbf{C}_{1,j}^p = \delta_0 \cdot w^0[j][0]$. The SSA output of neuron-(0) in layer-(1), as a function of its input features and filter weights can be computed as:

$$\begin{aligned}\nabla O^1[0] &= \sum_{j=0}^{N_0-1} (\nabla \mathbf{C}_{1,j}^p) \times w^1[0][j] = \sum_{j=0}^{N_0-1} (\delta_0 \cdot w^0[j][0]) \times w^1[0][j] \\ &= \delta_0 \cdot \sum_{j=0}^{N_0-1} (w^0[j][0]) \times w^1[0][j] \quad (12)\end{aligned}$$

Thus, independent of δ_0 , it results in exactly one equation, which is not sufficient to solve for N_0 unknowns. Going further, if we consider the multiple perturbation scenario:

$$\mathbf{C}^p(0, 1, \dots, (N_1 - 1)) = \begin{bmatrix} c_0 + \delta_0 & c_1 + \delta_1 & \dots & c_{(N_1-1)} + \delta_{(N_1-1)} \end{bmatrix}$$

it is possible to generate N_1 linearly-independent equations, but this system is also under-determined if $N_1 < N_0$.

2) *Solution*: This challenge can be solved through *localization*, where we place exactly one neuron in layer-(0) in the linear region of operation. For example, if we are interested in extracting $w^1[0][j]$ (the $(j+1)^{th}$ weight in $\mathbf{W}^1[0]$, i.e., the filter corresponding to neuron-(0) in layer-(1)), we need to search for the *Q-point* that **does** activate only the $(j+1)^{th}$ neuron in layer-(0), and **does not** activate the remaining neurons in layer-(0). In that case, Equation 12 simplifies to:

$$\nabla O^1[0] = \delta_0 \cdot (w^0[j][0] \times w^1[0][j]) \quad (13)$$

and hence the unknown weight in the second layer neuron can be extracted using the known values as:

$$w^1[0][j] = \frac{\nabla O^1[0]}{(\delta_0 \cdot w^0[j][0])} \quad (14)$$

In general, any weight $w^1[m][j]$ (the weight at the $(j+1)^{th}$ position of $\mathbf{W}^1[m]$, i.e., the filter corresponding to the $(m+1)^{th}$ neuron in the second layer) can be extracted/attacked in three steps.

- 1) Finding *Q-point* that activates only the $(j+1)^{th}$ neuron in the first layer, and does not activate remaining neurons in the first layer.
- 2) Applying perturbation δ_k at any arbitrary location-(k) in the input \mathbf{I} to the neural network.
- 3) And subsequently observing $O^1[m]$, i.e., the activated output of $(m+1)^{th}$ neuron in the second layer.

and finally extracting the model weight using:

$$w^1[m][j] = \frac{\nabla O^1[m]}{(\delta_k \cdot w^0[j][k])} \quad (15)$$

Although this subsection discusses the second layer, this *localization* method to extract model weights is generic and can be extended to any layer of the FCNN, *iteratively*, by searching for a *Q-point* that activates only one neuron in the prior layer. Once the weights are extracted, the biases can be successfully extracted similar to the technique used in Equation 11.

E. Linear Constraint Satisfaction

The activation of exactly one neuron in the first layer corresponding to the weight-under-attack (WUA) in the second layer neuron is basically a *constraint satisfaction* problem. The neuron function is itself linear, hence the constraint corresponding to the WUA index is a linear constraint with $>$ *inequality* and constraints for the remaining weights are also linear with \leq *inequality*. Hence, all the constraints are linear. Thus, we can formulate the problem of selective neuron's activation as a linear constraint satisfaction (LCS) problem². In order to extract model

²Without SLC formulation, the results will not converge. A random search on *IBM BladeCenter[®] High-Performance Cluster* visited > 10 billion solutions occupying 1TB of disk-space without convergence after 1 week of execution.

parameters with high fidelity corresponding to cases-(1),(2),(3), we will need to generate SLC that activates only one neuron in the prior layer (layer-(0)), and can be described as:

$$\begin{aligned}\mathbf{W}^0[j] \cdot \mathbf{I}^\top &> -1 * b_j + \eta, & j = j_{WUA} \\ \mathbf{W}^0[j] \cdot \mathbf{I}^\top &\leq -1 * b_j + \alpha, & j \in [0, N_0 - 1] - \{j_{WUA}\}\end{aligned} \quad (16)$$

However, in order to extract model parameters corresponding to case-(4), since isolated activation leads to the oracle response being "0", we will need to generate SLC that activates multiple neurons in the prior layer and reformulate on-the-fly as:

$$\begin{aligned}\mathbf{W}^0[j] \cdot \mathbf{I}^\top &> -1 * b_j + \tau, & j = j_{WUA}, \tau \in [v, \kappa] \\ \mathbf{W}^0[j] \cdot \mathbf{I}^\top &> -1 * b_j + \eta, & j = j_{fp}, \eta \gg \kappa \\ \mathbf{W}^0[j] \cdot \mathbf{I}^\top &\leq -1 * b_j + \alpha, & j \in [0, N_0 - 1] - \{j_{fp}, j_{WUA}\}\end{aligned} \quad (17)$$

Here fp corresponds to the first positive weight in the layer-(1) neuron and $v < \kappa \ll \eta$ to ensure the effect of negative weight will not dominate that of the positive weight, and hence ensure ReLU activation. This analysis is generic and extends to any layer-(i). Thus, for each neuron in layer-(i), (a) we generate an SLC for selective neuron activation in the layer-($i-1$), (b) solve the SLC, (c) extract the solution and make oracle queries to capture outputs of layer-(i) through scan-chains and finally (d) apply SSA to extract the model parameters of layer-(i). These steps need automation to enable *model-extraction* of large FCNNs of arbitrary depths and sizes.

*Algorithm 1*³ shows the proposed automation to extract $\hat{\theta}_{(i)}$ (model-subset for layer-(i)), given $\hat{\theta}_{(i-1)}$ (model-subset for already extracted layer-($i-1$)). Apart from $\hat{\theta}_{(i-1)}$, the algorithm also accepts other user-tunable parameters δ (the small-signal), η (target for cases-(1),(2),(3)), $[v, \kappa]$ (target range for case-(4)), α (to prevent numerical errors due to the neurons being at the boundary between deactivation and activation). During our experiments, we have also encountered *hard cases* (case-(5)) when the solver could not successfully solve some of the SLC instances and returns "Infeasible".

We note that the algorithm has different procedures for different cases. For cases-(1), (2), (3) (discussed in Section IV-C), isolated activation is used where we use an exclusive *Q-point* for higher fidelity. For case-(4), we are forced to use a common *Q-point*. To address case-(5), we have devised a strategy to relax constraints in Equations 16 and 17 for values of i corresponding to weights already successfully extracted. This is acceptable, because in Equation 12, we can substitute the extracted parameters and solve for the unknown parameters.

V. EXPERIMENTAL SETUP AND RESULTS

To evaluate and compare our results with earlier work using mathematical model extraction, we follow the approach of Carlini *et al.* [24]. We build the same neural network configurations with the same training process, apply our automated attack, and report the number of queries used and the maximum error rate ($\max|\theta - \hat{\theta}|$) of the extracted parameters. This will effectively extend their results table and compare our work with the three earlier attacks proposed at CRYPTO 2020 [24], USENIX 2020 [25], and ICML 2020 [26].

Table I presents the results of our attack and its comparison with the mathematical model extraction. The results clearly

³Simplified notation compared to Section IV-E.

Algorithm 1: Algorithm to attack target layer in FCNNs

```

Input:  $oracle, \hat{\theta}_{(i-1)}, \delta, \eta, \kappa, \nu$  ( $\eta \gg \kappa$ )
 $m := 0;$ 
while  $m < N_i$  do
   $j := 0;$ 
  while  $j < N_{(i-1)}$  do
    Generate constraints that activate only neuron-( $j$ ) in
    layer-( $i-1$ ) with  $\eta$  as target activation;
    Run constraint-solver for  $Q$ -point search;
    // Cases- (1), (2), (3)
    Apply  $SSA(oracle, Q, \delta)$  to extract  $w^i[m][j]$  weight;
    Mark index of first positive weight of neuron-( $m$ ) in layer-( $i$ )
    as  $fp$ ;
    // Coupled Constraints for case-(4)
    if extracted weight = 0.0 then
      Generate constraints that activate neuron-( $j$ ) with target
      activation  $\tau$  ( $\{v \leq \tau \leq \kappa\}$ ) and neuron-( $fp$ ) with target
      activation  $\eta$  (and none of the others) in layer-( $i-1$ );
      Run constraint-solver for  $Q$ -point search;
      Apply  $SSA(oracle, Q, \delta)$  to extract the negative weight.
    end
     $j := j + 1;$ 
  end
  // Iterative constraint relaxation for
  case-(5)
   $j := 0;$ 
  while  $j < N_{(i-1)}$  do
    if Constrained solution "Infeasible" then
      Relax constraints for already solved weights;
      Run constraint-solver for  $Q$ -point search;
      Apply Incremental extraction() along with
       $SSA(oracle, Q, \delta)$  to extract the weight.
    end
     $j := j + 1;$ 
  end
  Apply  $SSA$  to extract bias of neuron-( $m$ ) in layer-( $i$ );
   $m := m + 1;$ 
end
Result:  $\hat{\theta}_{(i)}$ 

```

TABLE I. Efficacy of our extraction attack, which is orders of magnitude more precise and query-efficient than prior work. Models denoted $a-b-c$ are fully connected neural networks with input dimension a , one hidden layer with b neurons, and c outputs. Entries denoted with a \dagger were unable to recover the network after ten attempts. Here, θ and $\hat{\theta}$ denote the original and extracted FCNN models respectively.

Architecture	Parameters	Approach	Queries	$\max \theta - \hat{\theta} $
784-32-1	25, 120	[25]	$2^{18.2}$	$2^{-1.7}$
		[24]	$2^{19.2}$	$2^{-30.2}$
		Ours	$2^{13.9}$	$2^{-53.3}$
784-128-1	100, 480	[25]	$2^{20.2}$	$2^{-1.8}$
		[24]	$2^{21.5}$	$2^{-29.4}$
		Ours	$2^{15.4}$	2^{-49}
10-10-10-1	210	[26]	2^{22}	2^{-12}
		[24]	2^{16}	2^{-36}
		Ours	$2^{8.8}$	$2^{-45.9}$
10-20-20-1	420	[26]	2^{25}	$\infty\dagger$
		[24]	$2^{17.1}$	2^{-37}
		Ours	$2^{10.2}$	$2^{-45.4}$
40-20-10-10-1	1, 110	[24]	$2^{17.8}$	$2^{-27.1}$
		Ours	$2^{11.1}$	$2^{-43.5}$
80-40-20-1	4, 020	[24]	$2^{18.5}$	$2^{-39.7}$
		Ours	2^{13}	$2^{-44.2}$

demonstrate the superiority of our attack across all the considered architectures. On average, our proposed attack is $2^{20.7} \times$, $2^{50.7} \times$ and $2^{33.9} \times$ more accurate and requires $2^{6.5} \times$, $2^{4.6} \times$ and $2^{14.2}$ less queries compared to [24], [25], and [26] respectively.

The results depend on the user-tunable parameters η, κ, ν ,

while the dependence on δ is observed to be small. We used the IBM® ILOG CPLEX *constraint-solver* for solving the SLCs. Algorithm 1 was implemented in C and the *constraint-solver* was invoked through system-calls. The results shown in Table I correspond to $\delta = 1.0$, $\alpha = -1.0$, $\eta = 10.0$, $\kappa = 2.0$ and $\nu = 1.0$ —these values are constant across all the FCNN architectures. The accuracy can be further improved by suitable choice of these parameters, which we leave for future investigation. Since ReLU also uses $\max(\cdot)$ operation, we assume ReLU activation-unit [23] can be re-used to compute the output layer, and hence the output of the last layer also leaks information. We have contacted authors of prior papers [24], [25] and confirmed that our assumption on the output layer is consistent with theirs.

Although it is an orthogonal threat, an interesting comparison of our attack can be made with the side-channel analysis, which seems to be more powerful as it can directly attack individual floating-point multiplications. It is reported that a few thousand queries are sufficient to extract the mantissa but direct comparison is hard because the attack was conducted on single-precision floating-point variables and the maximum error rate and queries are not given for the complete network [8].

VI. CONCLUSIONS AND FUTURE WORK

This paper has uncovered and quantified the effectiveness of a new threat vector for model stealing from ML hardware. The proposed attack abuses the debug ports (i.e., scan-chain access) that may be active after deployment to conduct in-field tests, and it observes intermediate states of the neural network inference for efficient model extraction. Although such attacks are well-known for cryptographic circuits, they have not been explored before for ML. We have revealed that ML has unique challenges due to the number of parameters, non-injective non-linear functions, and the floating-point arithmetic, which are not seen in cryptographic systems. But these challenges can be overcome with novel attack algorithms such as the ones we have proposed and with automation. The results have shown clear advantages of our proposed attack over cryptanalytic extraction both in terms of accuracy and the number of queries needed. This paper, therefore, calls for defenses against such attacks, e.g., by extending the ones built for cryptographic circuits.

To be consistent with earlier attacks and to evaluate NN achieving the best accuracy results, we used NN with floating-point weights and biases. An interesting extension of this work could be analyzing quantized networks.

REFERENCES

- [1] C. Li, "OpenAI's GPT-3 language model: A technical overview," 2020. [Online]. Available: <https://lambdalabs.com/blog/demystifying-gpt-3/>
- [2] F. Tramèr et al., "Stealing Machine Learning Models via Prediction APIs," in *USENIX Security Symposium*, 2016, pp. 601–618.
- [3] B. Wang and N. Z. Gong, "Stealing hyperparameters in machine learning," in *IEEE Symposium on Security and Privacy (SP)*, May 2018, pp. 36–52.
- [4] M. Juuti et al., "PRADA: Protecting Against DNN Model Stealing Attacks," in *IEEE European Symposium on Security and Privacy*, 2019, pp. 512–527.
- [5] D. Lowd and C. Meek, "Adversarial learning," in *ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, 2005, pp. 641–647.
- [6] N. Papernot et al., "Practical black-box attacks against machine learning," in *AsiaCCS*, 2017, pp. 506–519.
- [7] T. Gu et al., "BadNets: Evaluating Backdooring Attacks on Deep Neural Networks," *IEEE Access*, vol. 7, pp. 47 230–47 244, 2019.
- [8] L. Batina, S. Bhasin, D. Jap, and S. Picek, "CSI NN: Reverse Engineering of Neural Network Architectures Through Electromagnetic Side Channel," in *USENIX Security Symposium*, 2019, pp. 515–532.
- [9] L. Wei et al., "I Know What You See: Power Side-Channel Attack on Convolutional Neural Network Accelerators," in *ACSAC*, 2018.
- [10] Y. Xiang et al., "Open DNN Box by Power Side-Channel Attack," *IEEE Transactions on Circuits and Systems II: Express Briefs*, 2020.
- [11] A. Dubey et al., "MaskedNet: A Pathway for Secure Inference against Power Side-Channel Attacks," *HOST*, 2020. [Online]. Available: <https://arxiv.org/abs/1910.13063>

- [12] H. Yu et al., "DeepEM: Deep Neural Networks Model Recovery through EM Side-Channel Information Leakage," in *HOST*, 2020.
- [13] A. Dubey et al., "BoMaNet: Boolean Masking of an Entire Neural Network," *ICCAD*, 2020. [Online]. Available: <https://arxiv.org/abs/2006.09532>
- [14] J. Breier, D. Jap, X. Hou, S. Bhasin, and Y. Liu, "SNIFF: Reverse Engineering of Neural Networks with Fault Attacks," *CoRR*, vol. abs/2002.11021, 2020.
- [15] "Texas Instruments: IEEE Std. 1149.1 (JTAG) Testability Primer," 2016. [Online]. Available: <https://www.ti.com/lit/an/ssya002c/ssya002c.pdf>
- [16] "Tessent Boundary Scan," 2019. [Online]. Available: <https://www.mentor.com/products/silicon-yield/products/boundary-scan>
- [17] "1149.1-2013 - IEEE Std. for Test Access Port and Boundary-Scan," 2013. [Online]. Available: https://standards.ieee.org/standard/1149_1-2013.html
- [18] B. Yang et al., "Scan based side channel attack on dedicated hardware implementations of Data Encryption Standard," in *ITC*, 2004, pp. 339–344.
- [19] M. Agrawal et al., "Scan Based Side Channel Attacks on Stream Ciphers and Their Counter-Measures," in *INDOCRYPT*, 2008, pp. 226–238.
- [20] K. Rosenfeld and R. Karri, "Attacks and Defenses for JTAG," *IEEE Design and Test of Computers*, vol. 27, no. 1, pp. 36–47, 2010.
- [21] R. Nara et al., "Scan-based attack against elliptic curve cryptosystems," in *ASP-DAC*, 2010, pp. 407–412.
- [22] J. D. Rolt et al., "A Novel Differential Scan Attack on Advanced DFT Structures," *ACM TODAES*, vol. 18, no. 4, pp. 36–47, 2013.
- [23] Y. Chen et al., "Eyeriss: An Energy-Efficient Reconfigurable Accelerator for Deep Convolutional Neural Networks," *IEEE JSSC*, no. 1, pp. 127–138, 2017.
- [24] N. Carlini et al., "Cryptanalytic Extraction of Neural Network Models," in *CRYPTO*, vol. 12172, 2020, pp. 189–218.
- [25] M. Jagielski et al., "High accuracy and high fidelity extraction of neural networks," in *USENIX Security Symposium*, 2020, pp. 1345–1362.
- [26] D. Rolnick and K. P. Körding, "Identifying Weights and Architectures of Unknown ReLU Networks," *ICML*, vol. abs/1910.00744, 2020.
- [27] R. Torrance and D. James, "The state-of-the-art in semiconductor reverse engineering," in *DAC*, 2011, pp. 333–338.
- [28] V. Nair et al., "Rectified Linear Units Improve Restricted Boltzmann Machines," in *ICML*. Omnipress, 2010, pp. 807–814.
- [29] S. E. Laux, "Techniques for small-signal analysis of semiconductor devices," *IEEE Transactions on Electron Devices*, vol. 32, no. 10, pp. 2028–2037, 1985.