

A Hard-Label Cryptanalytic Extraction of Non-Fully Connected Deep Neural Networks using Side-Channel Attacks

Benoît Coqueret

Thales ITSEF

Toulouse, France

University of Rennes, INRIA, IRISA

Rennes, France

Email: benoit.coqueret@thalesgroup.com

Mathieu Carbone

Thales ITSEF

Toulouse, France

Email: mathieu.carbone@thalesgroup.com

Olivier Sentieys

University of Rennes, INRIA, IRISA

Rennes, France

Email: olivier.sentieys@inria.fr

Gabriel Zaid

Thales ITSEF

Toulouse, France

Email: gabriel.zaid@thalesgroup.com

Abstract—During the past decade, Deep Neural Networks (DNNs) proved their value on a large variety of subjects. However despite their high value and public accessibility, the protection of the intellectual property of DNNs is still an issue and an emerging research field. Recent works have successfully extracted fully-connected DNNs using cryptanalytic methods in hard-label settings, proving that it was possible to copy a DNN with high fidelity, *i.e.*, high similitude in the output predictions. However, the current cryptanalytic attacks cannot target complex, *i.e.*, not fully connected, DNNs and are limited to special cases of neurons present in deep networks. In this work, we introduce a new end-to-end attack framework designed for model extraction of embedded DNNs with high fidelity. We describe a new black-box side-channel attack which splits the DNN in several linear parts for which we can perform cryptanalytic extraction and retrieve the weights in hard-label settings. With this method, we are able to adapt cryptanalytic extraction, for the first time, to non-fully connected DNNs, while maintaining a high fidelity. We validate our contributions by targeting several architectures implemented on a micro-controller unit, including a Multi-Layer Perceptron (MLP) of 1.7 million parameters and a shortened MobileNetv1. Our framework successfully extracts all of these DNNs with high fidelity (88.4% for the MobileNetv1 and 93.2% for the MLP). Furthermore, we use the stolen model to generate adversarial examples and achieve close to white-box performance on the victim’s model (95.8% and 96.7% transfer rate).

1. Introduction

During the last decade, the number of tasks for which Deep Neural Networks (DNNs) have proven their effectiveness has steadily increased, leading to the widespread adoption of these algorithms in a large variety of fields.

From computer vision to text translation and images generation, DNNs are everywhere now, and the best models have become valuable intellectual property (IP). In the meantime, the parallel effort from hardware designers have made possible DNNs’ deployment on edge device. However, due to their high value, the IP of the deployed models must be protected against new attacks caused by the embedded context. Since the publication of a first side-channel attack against the IP of an embedded DNN [2], the number of physical-based attacks against DNNs has greatly increased. Several methodologies using side-channel attacks with the objective of the extraction of the DNN’s hyperparameters have been proposed [11], [18]. DNN parameters have also been targeted by physical attacks, via side-channel [19], [32], or through fault injection [13], [27]. These types of attack, targeting the parameters of the model, try to copy the targeted model and perform a model extraction attack.

Model extraction is not only a threat to embedded DNNs but to any deployed DNNs. There is therefore a large variety of methods, and even objectives, for these attacks. We can characterize the two main types of objectives or adversarial goals for model extraction with the terminology introduced in [17]: *accuracy-based model extraction* and *fidelity-based model extraction*. The first aims at gaining access to a substitute model with good performance on the task of the targeted model without having to perform the whole training process. The second has for purpose to clone the targeted model to acquire a copy as close as possible to the original. The cloned model can then be used to gain information on the victim’s DNN and potentially mount more powerful attacks against it. In this study, we will consider only fidelity-based model extraction.

Jagielski *et al.* [17] were the first to propose a functional framework for fidelity-based model extraction of 1-layer neural network (NN) using the ReLU (Rectified Linear Unit)

function. They exploited the gradient of the DNN to gain access to what they defined as *critical* points. Such points correspond to inputs where the activation value of one specific neuron is null. These points are visible in DNNs based on the ReLU activation function, since they correspond to points where discontinuities occur in the gradient of the DNN. Using this knowledge, they were able to extract a 1-layer fully-connected NN using the least-square (LSTSQ) algorithm. Due to the similarities with attacks targeting crypto-system, *i.e.*, analysis of a large number of input-output pairs to gain information on a secret value, this method was later characterized as *cryptanalytic extraction* of DNNs. Both works in [6] and [28] successfully extended this result to deeper architecture, with the method proposed in [6] still achieving the highest fidelity today. One key limitation in the method proposed in [6] was the extraction of the neurons’ sign, for which they used an exhaustive search. This is a critical information, as assigning the wrong sign will deactivate the neuron when it should be activated, and inversely. Shamir *et al.* argued that their method would not scale well with larger architecture [4]. So they used crafted perturbations activating specific neurons to infer this information [4]. While very effective, all of these methods rely on the fact that the output of the DNN is composed of the full confidence score vector which allows gradient estimation. This corresponds to an ideal case for the attacker, and researches have aimed at removing this assumption in order to move on to more realistic scenarios. Recent work by Chen *et al.* [9] demonstrated that it was possible to extract DNNs in hard-label settings for small networks. This was further proved by the work of Carlini *et al.* [5], in which they successfully extracted a four hidden-layer DNN using only the hard-label. However, while these works no longer require the confidence score, they are all designed solely for fully-connected DNNs, excluding any target composed of a non-fully connected layer, such as a pooling layer.

In parallel, several methods aiming for fidelity-based model extraction have been proposed using physical attacks. As mentioned before, the most notable examples of such attacks were presented in [13], [27]. They used fault injection to determine the value of a subset of the weights’ bits in the DNN and constraint learning for the rest of the weights. Contrarily to the methods previously mentioned, these are not restricted to fully connected DNNs and, as such, were successful in extracting various complex architectures, *e.g.*, ResNet-34 or VGG-11 [27]. However, these methods are either limited to DRAM platforms performing DNN inference, *i.e.* Machine-Learning-as-a-Service (MLaaS) platforms, [27] or require access to an open device, on which the attacker has full control (*i.e.* white-box settings), to find the memory localization of the weights’ bits [13].

We summarize in Table 1 the results and the threat models from state-of-the-art (SOTA) frameworks performing *fidelity-based model extraction* using side-channel, fault injection or cryptanalytic methods. To the best of our knowledge, there is no method for extracting complex architectures, *i.e.*, not restricted to fully-connected layers, without requiring an access to the confidence scores and an open

device.

Contributions: In this paper, we present an efficient framework, combining black-box side-channel attacks and cryptanalytic-based extraction of DNNs in a hard-label setting. We propose a new method to acquire critical points using a side-channel attack instead of the output of the DNN. This allows the attack to be performed in a hard-label settings, and offers higher precision than previous methodologies in extracting weights, resulting in higher fidelity. Additionally, our framework is not impacted by the output’s data format, and we are the first to propose results of cryptanalytic extraction with both 32- and 64-bit data. Furthermore, we introduce an alternative methodology to determine the sign of each neuron requiring only one hypothesis by layer. Finally, the main advantage of our method over other cryptanalytic extraction frameworks is the ability of the side-channel attacks to subdivide the DNN at each activation functions. Using this result, we improve the SOTA by targeting complex DNNs composed of non-fully connected layers, such as depth-wise separable convolution or pooling layers, in a MobileNetv1 architecture. We summarize the major contributions of our work as follows:

- We present a new black-box side-channel attack on a constant-time implementation of the ReLU function introduced as a SOTA countermeasure.
- We propose a new gradient-free extraction method improving both the precision on the weights’ extraction and the robustness to special cases of neurons.
- We introduce a new method to infer the sign of the neurons without access to the confidence scores and requiring only the testing of one hypothesis by layer.
- We build an end-to-end framework capable of extracting DNNs with high fidelity and not limited to fully-connected layers.
- To prove the practicability of our contributions, we validate them through comparison against other SOTA frameworks, and by targeting a shortened version of MobileNetv1 embedded on an STM32F767ZI using the X-Cube-AI framework.

We provide an implementation of our code at <https://github.com/X>.

2. Background

2.1. Notations

Let calligraphic letters \mathcal{X} denote sets, the corresponding capital letters X (resp. bold capital letters) denote random variables (resp. random vectors \mathbf{T}), and the lowercase x (resp. \mathbf{t}) denote their realizations. We will use \mathbf{T}_i to describe the i -th element of a vector \mathbf{T} . Throughout this paper, the function modeled by a DNN is denoted as $f : \mathcal{X} \rightarrow |\mathcal{Y}|$, which characterizes its ability to classify a data $X \in \mathcal{X}$, *e.g.*, an image, over a set of $|\mathcal{Y}|$ classes. A DNN characterized by the function f and the set of weight θ is denoted f_θ , and we denote $f_{\hat{\theta}}$ the model extracted from f_θ . The probability of observing an event X is denoted by $\Pr[X]$. Finally, we denote $\mathcal{D}_{\mathcal{X}}$, the distribution over the set of data \mathcal{X} .

TABLE 1: Overview of the state-of-the-art of fidelity-based model extraction attacks.

Approach	Attack type	Full extraction (weight + bias)	Hard-label setting	Not restricted to fully connected DNN	Random queries	DNN's datatype tested	Targeted architecture (Most complex)
ICML'20 [28]	Cryptanalytic	✓	✗	✗	✓	64-bit float	MLP 10-20-20-1
Crypto'20 [6]	Cryptanalytic	✓	✗	✗	✓	64-bit float	MLP 40-20-10-10-1
AC'24 [9]	Cryptanalytic	✓	✓	✗	✓	64-bit float	MLP 1024-2-2-1
Preprint [5]	Cryptanalytic	✓	✓	✗	✓	64-bit float	MLP 3072-256×3-64-10
USENIX'19 [2]	Side-Channel	✗	✓	✓	✓	32-bit float	MLP 784-200×4-10
ICCAD'23 [33]	Side-Channel	✓	✗	✗	✓	64-bit float	LeNet5
IEEE S&P'22 [27]	Fault Injection	✓	✓	✓	✗	8-bit data	ResNet34 and VGG11
ESORICS'23 [13]	Fault Injection	✓	✗	✓	✗	8-bit data	CNN 3×(Conv + Pooling + ReLU)-Linear
This work	Cryptanalytic and side-channel	✓	✓	✓	✓	32- and 64-bit float	Shortened MobileNetv1 MLP 3072-256×3-64-10

Attacks not in hard-label settings suppose an access to the confidence scores. Contribution not using random queries use a subset of the training or a testing dataset to perform part of their attacks.

2.2. Fidelity-based model extraction

Model extraction attacks target the confidentiality of a deployed model, by trying to obtain a copy of the victim’s model. The most common goal, is to use the copy to gain the target’s benefits, *e.g.*, commercial value, performance on specific tasks, *etc.*, without having to train a model. Another possibility is to use the stolen model to gain information on the target and mount higher-level attacks against it, *e.g.*, generation of adversarial examples in a white-box scenario [7], [29] or membership inference attacks [10]. This diversity in the attack scenarios leads to several adversarial goals in a model extraction. Following the notations introduced in [17], we consider two main goals, namely *task accuracy* and *fidelity*.

The first one has for purpose to extract $f_{\hat{\theta}}$, such as, for the true task distribution $\mathcal{D}_{\mathcal{X} \times \mathcal{Y}}$ over the sets of input \mathcal{X} and label \mathcal{Y} , $f_{\hat{\theta}}$ maximizes $\Pr_{\mathcal{X}, \mathcal{Y}}[\arg\max(f_{\hat{\theta}}(X) = Y)]$. In practice, this can be achieved through learning-based methods and optimisation. However, this approach induces variability in the optimisation problem which can lead to very different solutions (*i.e.*, models) caused by the convergence to different local minima. This is why this method offers no guaranty from a *fidelity* point of view, and is used for the goal of *task accuracy*.

The second adversarial goal consists in maximizing the similarity between the stolen and the original models in their predictions. The limits being *functionally equivalent* extraction which is defined by:

Definition 1 (Functional equivalence [17]). Two models f_{θ} and f_{γ} achieve functional equivalence on \mathcal{X} if for any distribution of input $\mathcal{D}_{\mathcal{X}}$, the following equality holds:

$$\forall X \in \mathcal{X}, f_{\theta}(X) = f_{\gamma}(X).$$

This is the strongest attack possible as it leads to an exact copy of the targeted DNN. In this paper, we only consider fidelity-based model extraction. To evaluate the success of such extraction, we use a relaxed definition of functional equivalence:

Definition 2 ((ϵ, δ)-functional equivalence [6]). Two models f_{θ} and g_{γ} are (ϵ, δ)-functionally equivalent on \mathcal{X} if:

$$\Pr_{X \in \mathcal{X}}[|f_{\theta}(X) - g_{\gamma}(X)| \leq \epsilon] \geq 1 - \delta.$$

The advantages of fidelity-based model extraction have stimulated researches on this subject. In particular physical-based attacks and cryptanalytic approaches have been proposed in recent years. Side channel attacks [2], [14], [19], [33] and fault-based attacks [13], [27] have been used to target all the weights of a DNN individually [2], [14], [19], or to extract partial information on the weights in order to train a substitute network using this information as a constraint [13], [27]. These methods achieve fidelity-based extractions but show limitations either in the attack complexity, in the architecture of the targeted DNN, or in the threat model, *e.g.*, exclusive to certain platforms, white-box settings such as access to the logits or to an open copy of the targeted device. On the other hand, attacks based on cryptanalysis have proved effective against shallow networks [17], [26], and were then successfully adapted to DNNs [6], [28]. In particular, Carlini *et al.* [6] were able to steal a DNN with three hidden layers using $2^{17.8}$ queries and extract the 1,110 parameters with a maximum absolute error between the weights θ and the stolen weights $\hat{\theta}$ of $2^{-27.1}$. These methods will be briefly introduced in the following section.

2.3. Cryptanalytic extraction methodology

Both methods in [6], [28] extract the DNN’s weight by targeting iteratively each neuron, monitoring its state, in order to find its critical points. The critical points and the state of a neuron η can be defined in the following way:

Definition 3 (Critical point and state of a neuron [6]).

Given an input $X \in \mathcal{X}$, let $V(\eta; X)$ be the function characterizing the input of the neuron η before applying the activation function. Then X is said to be a critical point to the neuron η if $V(\eta; X) = 0$. If $V(\eta; X) > 0$ (resp. $V(\eta; X) < 0$) then η is said to be active (resp. inactive).

The induced-hyperplane of a neuron can be defined accordingly:

Definition 4 (Neuron-induced Hyperplane). The hyperplane of a neuron η corresponds to the set $\mathcal{H}_\eta \subseteq \mathcal{X}$ where $\forall X \in \mathcal{H}_\eta, V(\eta; X) = 0$.

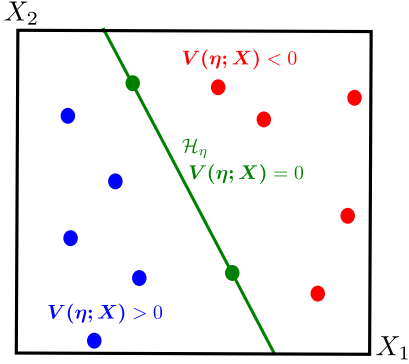


Figure 1: Neuron-induced hyperplane.

An example of critical points, neuron states and hyperplane is provided in Figure 1. By finding enough critical points, it is possible to determine the equation of the neuron-induced hyperplane. This can then be used to find the weights' values with high precision, since by definition, the equation of the hyperplane follows:

$$\sum_{i=1}^n \theta_{\eta,i} X_i + \beta_\eta = 0, \mathbf{X} \in \mathbb{R}^n \quad (1)$$

with θ_η (resp. β_η) the weight (resp. the bias) vector associated with the neuron η . It is important to notice that for a neuron η with n parameters, the induced-hyperplane is an $(n - 1)$ -dimensional piecewise-linear surface [12]. This causes that the weight vector can only be extracted up to a scaling factor. In practice, we set the first weight $\theta_{\eta,0}$ of θ_η as the scaling factor and from Equation 1 we obtain:

$$\mathbf{X}_0 = \sum_{i=1}^n \frac{-\theta_{\eta,i}}{\theta_{\eta,0}} X_i - \frac{\beta_\eta}{\theta_{\eta,0}}, \theta_{\eta,0} \neq 0 \quad (2)$$

After the extraction of $n - 1$ critical points, it is possible to retrieve the neuron's weights up to a scaling factor. This process is defined in [6] as signature search. The neuron's signature can be defined in the following manner:

Definition 5 (Neuron signature [4]). The signature of a neuron η corresponds to a vector $\mathbf{S}_\eta \in \mathbb{R}^n$, for which there exists a scalar $\alpha \in \mathbb{R}$ such that $\theta_\eta = \alpha * \mathbf{S}_\eta$.

In the previous example, the vector $\mathbf{S}_\eta = (1, \frac{\theta_{\eta,1}}{\theta_{\eta,0}}, \dots, \frac{\theta_{\eta,n}}{\theta_{\eta,0}})$ is a signature of the neuron η and $\theta_{\eta,0}$ is the scaling factor. By considering the target model f_θ as an oracle, the search for neuron's signature can be performed directly using its input-output pairs. Indeed, if the target model is composed of piecewise linear activation functions, then f_θ is also piecewise linear and discontinuities can be observed in the gradient of f_θ . These discontinuities correspond to each change of state by a neuron in the DNN. However performing the signature's search using these discontinuities requires high

assumptions. Firstly, the target model's output needs to be either the confidence scores or the logits, to allow the attacker to estimate the gradient through finite difference. Secondly, the output needs to be returned in a high-precision format to get the most precise estimation. This method was introduced in [17] and adapted to DNNs in [6], [28] to extract ReLU-based network running double-precision, 64-bit floating-point data.

Targeting deep ReLU networks is motivated by both the popularity of the ReLU function as an activation function and its properties. In particular, the discontinuity in the gradient of the ReLU function is localized in 0, which implies that by monitoring the gradient of the targeted model, critical points for each neuron can be found. Furthermore, the ReLU function is equivariant under positive multiplication, *i.e.*, $\forall x \in \mathbb{R}$ and $\forall c \in \mathbb{R}^+$, $ReLU(c \times x) = c \times ReLU(x)$. This means that the value of the scaling factor α of the neuron's signature is not important, only its sign. If the sign of the input of the ReLU function is respected, then the activation value of each neuron can later be re-scaled during the extraction of the neurons' signature in the following layers. Therefore the model extraction attack can be seen as an iterative process over the neurons based on two steps:

- 1) Extraction of the neuron's signature.
- 2) Extraction of the scaling factor to ensure that the neuron's state remains the same.

While neuron signatures can be revealed through the discontinuities of the gradient, the sign of the scaling factor is not accessible directly by an attacker. Carlini *et al.* [6] solve this by testing all of the 2^m possibilities with m the number of neurons in the layer. As this becomes rapidly impractical for large DNNs, Canales-Martínez *et al.* [4] propose another method based on activating specific neurons through crafted inputs to extract the sign of each neuron. Their method successfully extracts the sign of each neuron within a model with 8 hidden layers each with 256 neurons, for which the signatures were previously extracted.

Although, these methods have successfully performed complete extraction of Deep-ReLU networks up to 3 hidden layers and 100,480 parameters [6], deeper extraction has not been achieved. Furthermore, access to the confidence scores or the logits in a high-precision format is essential to the success of the extraction, which will make the attack impractical in most contexts. For example, in embedded systems, the hardware constraints can impose limitations on the precision or on the access to the logit values that could prevent the attack. However, the embedded-system context increases the surface of attack and allows an attacker to consider side-channel attacks to improve fidelity-based model extractions.

2.4. Side-channel attacks

Historically, side-channel analysis (SCA) is a class of cryptographic attack in which an attacker tries to exploit the vulnerabilities of the implementation of a real-word cryptosystem for key recovery by analyzing its physical characteristics *via* side-channel traces, like power consumption or

electromagnetic (EM) emissions. During the execution of an algorithm embedded into a crypto-system, side-channel traces record the intermediate variable (*e.g.*, secret key) being processed. To discriminate the correct key hypothesis from the incorrect ones, an attacker uses a statistical distinguisher to extract the targeted variable from a large number of traces.

Definition 6 (Statistical distinguisher). A statistical distinguisher, denoted $d : \mathcal{T} \rightarrow \mathcal{K}$, is a statistical function identifying the dependence between a set of physical traces \mathcal{T} and a targeted variable in \mathcal{K} .

In [2], Batina *et al.* were among the first to transpose this attack to the DNN paradigm by targeting the architecture and the parameters of the model with two statistical distinguishers via a Simple Power Attack (SPA) and a Correlation Electromagnetic Attack (CEMA). Since, CEMA has been extended to target the weights of DNN on microcontrollers (MCUs) [19] as well as on GPUs [14]. Several other methodologies commonly applied in side-channel attacks on crypto-system have been extended to target DNNs. In particular, deep learning-based side-channel attacks [3], [22] have been successfully adapted for extracting the hyperparameters of DNNs. Gao *et al.* [11] trained a set of meta-models on physical traces to extract the architecture of the targeted models with a high precision. Previously, Maia *et al.* [23] successfully reconstructed the architecture of ResNet and VGG models running on a GPU using the physical traces. All these attacks use supervised learning to construct the statistical distinguisher. This imposes to the attacker an access to an open device that provides the true label before learning the mapping. This assumption is very restrictive, leading researchers to focus on more permissive threat models using unsupervised learning methods. In particular, well-researched unsupervised cluster classification algorithms, such as *k-means* clustering [1], can be used as statistical distinguisher to find partitions effectively without any manual methods or prior profiling. While we argue that this method could improve fidelity-based model extraction, to the best of our knowledge, it has never been applied to the extraction of DNN weights.

3. Limitations in fidelity-based extraction

Although the methods presented in the previous section demonstrate strong results, fidelity-based model extraction remains a difficult task, especially in the context of embedded systems. In Section 3.1, we provide further details on the limitations of these methodologies in this specific context. Then, in Section 3.2, we describe the threat model that is considered in this paper. Finally, in Section 3.3, we present the overall flow of our fidelity-based extraction framework.

3.1. Fidelity-based extraction issues

Even if an ideal scenario is assumed for the attacker, *i.e.*, access to the complete confidence scores is given, cryptanalytic extraction using the previously introduced method is

still challenging for several reasons. First of all, the information in the input-output pairs does not allow identification of the neuron for which the critical point is found. Therefore, an attacker cannot validate if consecutive critical points found are related to the same neuron. Carlini *et al.* [6] solve this issue by sampling a large number of critical points with respect to the number of neurons. This statistically ensures that each neuron has enough critical points to infer the hyperplane required to find the corresponding signature, but leads to the search of unnecessary critical points. This might cause the attack to be too complex for deeper architectures regarding the large number of required queries. Another issue, reported in [4], is the special cases of neurons that almost never change state. For these neurons, the number of queries needed for the search of critical points can exceed the mean number of required queries for a classical neuron by a factor of 10. Even if these special cases of neuron are highly difficult to threat in practice, no investigations are provided in the SOTA to deal with such issue.

Additional practical limitations should be considered in the embedded context, as hardware constraints might force the neural network to use low-precision format of data for the weights, the activations and the input-output pair. This makes gradient estimations noisier, which has an impact on the search for critical points. Furthermore, in the case of classification task, if the output of the model f_θ is only the hard-label, *i.e.*, $\hat{y} = \operatorname{argmax}(f_\theta(X))$, then the attack scenarios designed in [6], [17], [28] are no longer applicable. While the logits or the confidence scores might be given alongside the prediction in a MLaaS context, for embedded systems, it is frequent that the output is restricted only to the class with the highest probability. For example, FINN [31] is a popular framework for the deployment of low-precision neural network on FPGA, and by default only the top-1 prediction is returned as output.

Finally, even if the application context of the neural network imposes that the output includes the confidence scores, several countermeasures can easily be implemented to increase the cost of the cryptanalytic extraction attack. Indeed, confidence scores are already used to generate adversarial examples in a black-box scenario [8]. Therefore, countermeasures designed to limit adversarial example generation, such as rounding or adding Gaussian noise to the confidence scores [16], [30], would also be efficient against extraction frameworks based on cryptanalysis such as [6], [17], [28].

All of these practical issues have driven recent researchers to find an alternative not requiring the gradient to find the critical points. Carlini *et al.* [5] propose to use the decision boundary to find dual points, *i.e.*, points that are both critical and on the decision boundary, and extract the trained parameters (*i.e.*, the weights and biases). In their work, they successfully extracted a four-layer DNN trained on CIFAR-10 using only the hard-label provided by the targeted model. However, their method is limited to fully-connected networks, and, as such, cannot be extended to any architecture using a layer that is not fully connected, *e.g.*, an average pooling layer. Furthermore, while they report

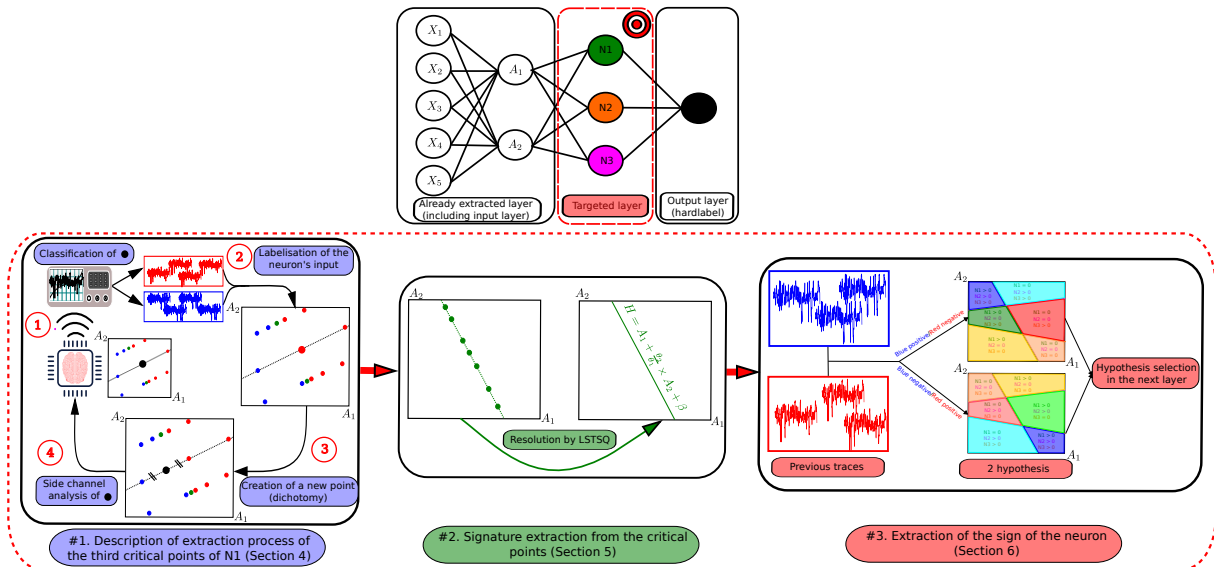


Figure 2: Overview of our attack framework. In Stage 1, the critical points of the neuron N1 are extracted using side-channel leakages. During Stage 2, the critical points, already extracted, are used to infer the signature of N1. Finally, the sign of the neuron is extracted in Stage 3.

the presence of special case neurons, they do not provide an alternative for their extraction, leading to an increase in the number of queries [5].

3.2. Threat model

In this paper, we consider an attacker, aiming for extraction of the weights of a DNN implemented in an embedded system (*e.g.*, an MCU). The purpose of this attack is to achieve functional equivalence between the victim’s model and the stolen copy (see Definition 1).

We assume that the attacker has a physical access to the device on which the targeted DNN is running. This allows recording of the physical leakages, *e.g.*, EM traces, during the processing of the input data by the DNN. However, contrary to [13], [33], we restrict ourselves to the case where the attacker does not have the possibility to profile the embedded device. As a consequence, it is not possible to train a statistical distinguisher on labeled physical traces. From the physical attack point of view, this setting corresponds to a black-box scenario.

Additionally, similarly to [17], [28] and [6], we only consider Deep-ReLU networks as a target in this threat model, and we assume that the exact architecture of the targeted model is known. However, we do not restrict ourselves to fully-connected layer. Finally, we suppose unrestricted query access to the embedded model, and so, the DNN can be seen as an oracle by the attacker. However, in this threat model, we examine the case where the output of the targeted DNN is restricted to the hard-label and not a vector composed of the class-probabilities or the logits.

3.3. Fidelity-based extraction assisted by side-channel attacks

In this work, we introduce a new attack framework against DNNs, using side-channel information to assist cryptanalytic extraction. It is an iterative process over all neurons of all layers. Our framework begins by targeting the input layer before moving to the first hidden layer, and so on. The extraction of each neuron follows the following three stages, as represented in Figure 2:

- 1) During the first step of our framework, we use side-channel information to infer the targeted neuron’s state for each input, and binary search to find the critical points. Instead of the gradient, this method exploits side-channel information through the use of a black-box statistical distinguisher (see Definition 6) to estimate the state of the neuron. Therefore, this stage is free from any of the restrictions mentioned in Section 3.1. This stage will be detailed in Section 4.
- 2) The second stage consists in retrieving the optimal equation describing the hyperplane induced by the targeted neuron. To do that, a system of equations is constructed based on the critical points extracted in the previous stage. Then, we use the least square algorithm (LSTSQ) on this system of equations to find the optimal solution, *i.e.*, the hyperplane’s equation. Further explanations on this part will be provided in Section 5.
- 3) In the final stage of our framework, the sign of each neuron is extracted through hypothesis selection, similarly to [6]. However, we propose to take leverage the physical traces collected during Stage 1. This reduces the number of hypothesis from 2^m , with m the number

of neurons in the layer, to a single hypothesis for the whole layer. This method will be presented in Section 6.

4. Stage 1: Search of Critical Points

Without any access to the confidence scores, the estimation of the gradient is not possible, and the methodologies to search for the critical points proposed in [6], [17], [28] are no longer applicable. Furthermore, the search for the dual points as introduced in [5], while effective in a hard-label settings, is only possible on fully-connected DNNs. This highlights the need for another method for the critical point search in hard-label settings for non-fully connected DNNs. In this section, we propose a new method, corresponding to the first stage of our framework (see Figure 2), to address these limitations. First of all, in Section 4.1, we present how we exploit physical leakages to extract the targeted neuron’s state without access to the confidence scores. Then, in Section 4.2, we describe the algorithm that we use to find the critical points associated with the neuron based on the previously determined neuron’s state.

4.1. Neuron’s state identification using side-channel information

In this study, we only consider Deep-ReLU networks, meaning that we restrict ourselves to DNNs where the activation function is always the ReLU function $\sigma(x) = \max(0, x)$. This is a common assumption, as the ReLU function is one of the most popular choice for activation function in DNNs. The key of the first step of our framework is based on constructing a statistical distinguisher mapping a side-channel trace to one of the two possible states of a neuron. The purpose of such distinguisher is to solve a classification problem and discriminates active from inactive neurons, based on physical traces characterizing their process. Then, once this separation is constructed, it is possible to infer the state of a given neuron for one input. Consequently, for two different inputs, using the statistical distinguisher allows the attacker to determine if they correspond to the same state for a given neuron, *i.e.*, if they are on the same side of the induced-hyperplane. It is important to notice that such distinguisher does not need to provide a mapping that returns the true state, *i.e.*, active or inactive. As we are only interested in the similarity in the neuron’s state, only the separation learned by the statistical distinguisher is required.

The activation function of DNNs has already been targeted using a statistical distinguisher exploiting a difference of timing in [2], [24], leading to the identification of the used activation function and even the weights’ mantissa. To prevent this issue, Maji *et al.* [24] proposed a constant-time implementation of the ReLU function based on a mask derived from the sign of the processed value. It is nevertheless worth mentioning, that non-constant-time implementations can still be found in popular open-source frameworks such as NNOM [21].

While effective, both previously mentioned approaches used supervised learning to find the best statistical distinguisher. This strategy requires access to an open device which is a very restrictive assumption. In this work, we consider a more realistic threat model in which an open device is not accessible to the attacker. Therefore, the supervised learning strategy is impossible.

To achieve this, we propose another kind of statistical distinguisher based on an unsupervised *k-means* clustering algorithm [1]. The idea is to first generate enough traces with random inputs to build a dataset in which both classes (*i.e.*, active and inactive neuron) are evenly represented. Then, we use the *k-means* algorithm to identify two clusters corresponding to each of the states. As mentioned previously, it does not matter to know which cluster corresponds to the active or inactive state, since that information is not needed to find a critical point. As the inputs are randomly generated, these two clusters are then used as reference to identify the state of a neuron. Finally, given the physical trace of a new neuron, we use its Euclidean distance to the two clusters’ centroids as our statistical distinguisher function to infer the neuron’s state.

4.2. Binary search of critical points

Using physical leakages, we are able to infer the state of a neuron through clustering. Even without the gradient, this approach, combined with binary search, is successful in finding the critical points. Indeed, starting from two random points, providing that they correspond to the two different states for the targeted neuron, an attacker can search for critical points by dichotomy, as represented in Figure 2. A high-level description of our methodology is presented in Algorithm 1.

Due to numerical approximations, finding a point on the hyperplane can be challenging. Instead, we use an hyperparameter Δ to place a critical point within a determined distance to the hyperplane. We use the *k-means* statistical distinguisher to infer the state of the neuron at each step of our binary search, and determine which point to update to get closer to the hyperplane. Since we are also using binary search to retrieve the critical points, the complexity of this phase is the same as the one in [17] and will require $\lceil \log_2(\Delta) \rceil$ queries to place a point within a distance of Δ to the hyperplane.

It is important to notice that the vulnerability that leaks the neuron’s state comes from the DNN’s implementation itself and is therefore independent from the output provided by the targeted DNN, contrary to classical cryptanalytic methods [6], [17]. Another important limitation of the classical cryptanalytic methods is their inability to determine to which neuron a critical point corresponds. To circumvent this issue, the solution adopted in [6] consists in sampling a large number of critical points, *i.e.*, the number of total collected critical points n_c is established via the coupon collector argument $n_c \gg N \log(N)$ with N the number of total neurons in the network, to ensure that there is enough critical points for each neuron to extract the signature. This

Algorithm 1 Binary search of critical point

Require: an input $X \in \mathcal{X}$, a statistical distinguisher d , a trace generator $Q : \mathbb{N} \times \mathcal{X} \rightarrow \mathcal{T}$ related to the targeted neuron $\eta \in \mathbb{N}$, distance precision $\Delta \in \mathbb{R}$,

Ensure: $X \in \mathcal{H}_\eta$

$Y \leftarrow \text{random}(\mathcal{X})$

$T_X, T_Y \leftarrow Q(\eta, X), Q(\eta, Y)$

$V(\eta; X), V(\eta; Y) \leftarrow d(T_X), d(T_Y)$

while $V(\eta; X) \neq V(\eta; Y)$ **do**

$Y \leftarrow \text{random}(\mathcal{X})$

$T_Y \leftarrow Q(\eta, Y)$

$V(\eta; Y) \leftarrow d(T_Y)$

end while

while $\|X - Y\|_2 > \Delta$ **do**

$M \leftarrow \frac{X+Y}{2}$

$T_M \leftarrow Q(\eta, M)$

$V(\eta; M) \leftarrow d(T_M)$

if $V(\eta; M) = V(\eta; X)$ **then**

$X \leftarrow M$

else ▷ Implies that $V(\eta; M) = V(\eta; Y)$

$Y \leftarrow M$

end if

end while

return X

issue is not present in our contribution, as we use a statistical distinguisher to extract the targeted neuron’s state. We are therefore able to localize to which neuron corresponds the identified critical point. Once all the desired critical points are found, an attacker can use them to retrieve the neuron’s signature.

5. Stage 2: Signature Extraction

In this section, we describe how we use the collected critical points to infer the targeted neuron’s signature. This corresponds to Stage 2 of our framework depicted in Figure 2. Section 5.1 introduces the general method. Then, in Section 5.2, we propose the first solution to deal with special neurons introduced in [4], namely *always-on* and *always-off*, and for a new type of special neurons, *input-off* neurons.

5.1. Generic method

The idea is to infer the equation of the neuron-induced hyperplane (see Definition 4) by constructing a linear system based on Equation 1 with the critical points from Stage 1. Such system can then be used to extract the weights. For example, given a set of n critical points $\{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^n\}$ s.t. $\mathbf{X}^i \in \mathbb{R}^n$, the following system of equations can be expressed in order to extract the n weights of the targeted

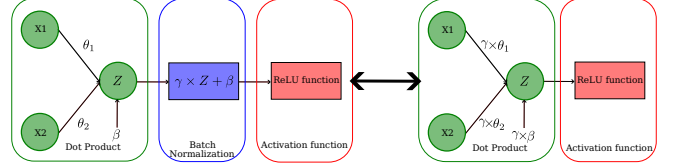


Figure 3: Signature extraction through batch normalization.

neuron η :

$$\begin{cases} \sum_{i=1}^n \frac{-\theta_{\eta,i}}{\theta_{\eta,0}} \mathbf{X}_i^1 - \frac{\beta_\eta}{\theta_{\eta,0}} = \mathbf{X}_0^1 \\ \vdots \\ \sum_{i=1}^n \frac{-\theta_{\eta,i}}{\theta_{\eta,0}} \mathbf{X}_i^n - \frac{\beta_\eta}{\theta_{\eta,0}} = \mathbf{X}_0^n \end{cases} \quad (3)$$

During the second stage of our framework, we find the equation of the hyperplane by solving the least square problem of this system of equations. We then extract the signature of the targeted neuron by arbitrarily fixing the scaling factor $\theta_{\eta,0}$ to ± 1 . The sign of the scaling factor determines the state of the neuron, and is therefore, a critical component to extract. While the critical point search does not require to label each cluster, we need to assign it to conduct the full extraction attack. To do so, the sign must be retrieved. This step will be detailed in Section 6.

It is important to notice that the proposed method is not impacted by the presence of linear transformations, *e.g.*, batch normalization, between the matrix multiplication and the activation function. Indeed, it is possible to fuse all these linear operations into one, and then extract an aggregate of the weights (see Figure 3).

This framework splits the neuron’s signature extraction process at each activation and is therefore the first framework not impacted by non-fully connected layers. However, one issue mentioned in SOTA’s frameworks which still impact our methodology, is the presence of special cases of neurons.

5.2. Special cases of neurons

From the previous section it can be observed that the whole signature’s extraction process is based on a system of equations constructed with the critical points of the targeted neuron η as represented in Figure 2. Some conditions must thus be respected. First of all, the number of critical points, *i.e.*, the number of equations, needs to be at least equal to the number of targeted parameters (weights, bias, *etc.*) in the neuron. And secondly, each targeted parameter must be expressed through a non-null component value in at least one critical point. In other words, the matrix describing the system of equations must be of full rank.

Unfortunately, there exist neurons for which it is difficult to meet these conditions, and which therefore require a very large number of queries to find a subset of critical points

satisfying them. Such neurons can be broadly classified in two types:

- Neurons that almost never change states, except on a very small subset of the input space \mathcal{X} , making it difficult to find critical points via random queries. Such neurons were already described in [4], but no solution was proposed to extract their signature. Neurons that are almost always active (resp. inactive) are referred as *always-on* (resp. *always-off*).
- Neurons for which one of the components of the critical points is almost always null. This specific type of neurons was not mentioned in [4], but cause issues for the signature extraction. Indeed, the number of queries needed to form a full-rank matrix for these neurons can exceed the average number of queries needed for classical neurons by a factor of 10. We designate such neurons as *input-off*.

For these types of special neurons, the signature extraction differs and a special process depending on the type of special neurons is performed. It is important to notice that, due to the information we extract from physical leakages, we are able to differentiate *always-off* from *always-on* neurons. This distinction is closely related to the sign extraction of the scaling factor using side-channel analysis (see Section 6).

Extraction of *always-off* neurons. The most simple case of specific neuron is the *always-off* neuron. Since its output value is almost always null and therefore independent from the input, we can directly assign its output value to 0 without having to extract the related weights.

Extraction of *always-on* neurons. For *always-on* neurons, the search for critical points can be extremely long and costly in queries. Furthermore, due to the sequential paradigm of our attack, applying our methodology on a layer without complete extraction of the previous ones is impossible. It is then crucial to extract all neurons at each layer. Since we know that the neuron is *always-on*, one solution is to consider the neuron’s activation function, *i.e.*, the ReLU function, as the identity and treat the neuron as a weighted skip connection between this layer and the next one. Let $\mathbf{X} \in \mathbb{R}^n$ be the input of a layer l composed of m neurons and η a *always-on* neuron in this layer l with $\boldsymbol{\theta} \in \mathbb{R}^n$ (resp. β_η) its weight vector (resp. bias). Let λ be a neuron in layer $l + 1$ and $\boldsymbol{\gamma} \in \mathbb{R}^m$ (resp. β_λ) its weight vector (resp. bias). If we denote \mathbf{A}_i the activation value of the neuron i in layer l corresponding to the input \mathbf{X} of layer l , the neuron η can be approximated as a skip connection to λ , and Equation 1 then becomes:

$$\sum_{\substack{i=1, \\ i \neq \lambda}}^m \gamma_i \mathbf{A}_i + \gamma_\eta \times \sum_{j=1}^n \boldsymbol{\theta}_j \mathbf{X}_j + \gamma_\lambda \beta_\eta + \beta_\lambda = 0 \quad (4)$$

A schematic description of the approximation process for *always-on* neurons is provided in Appendix A. As we can see, this equation remains linear with respect to the targeted neuron’s weights. Based on Equation 4, a new system of equations can be constructed for the targeted *always-on*

neuron. Using the same notations as in Equation 4, we can describe it in the following manner:

$$\begin{cases} -\sum_{i=1}^m \frac{\gamma_{\lambda,i}}{\gamma_{\lambda,0} \boldsymbol{\theta}_{\eta,0}} \mathbf{A}_i^1 - \sum_{j=1}^n \frac{-\boldsymbol{\theta}_{\eta,j}}{\boldsymbol{\theta}_{\eta,0}} \mathbf{X}_j^1 - \frac{\beta_{\eta\lambda}}{\boldsymbol{\theta}_{\eta,0}} = \mathbf{X}_0^1 \\ \vdots \\ -\sum_{i=1}^m \frac{\gamma_{\lambda,i}}{\gamma_{\lambda,0} \boldsymbol{\theta}_{\eta,0}} \mathbf{A}_i^z - \sum_{j=1}^n \frac{-\boldsymbol{\theta}_{\eta,j}}{\boldsymbol{\theta}_{\eta,0}} \mathbf{X}_j^z - \frac{\beta_{\eta\lambda}}{\boldsymbol{\theta}_{\eta,0}} = \mathbf{X}_0^z \end{cases} \quad (5)$$

with $\beta_{\eta\lambda}$ an aggregate of the biases of η and λ , and z the number of total equations (s.t. $z \geq n + m - 2$).

As we can see, the system of equations is very similar to the one described in Equation 3. Therefore, we use the same methodology to solve it. We can also observe that the signature of the neuron η appears in the extracted signature of λ . So, by solving the least square problem for a neuron in the next layer, it is possible to find the signature of an *always-on* neuron. The only exception is the bias which cannot be directly extracted. Instead, solving the system of equations returns a combination of the bias from η and λ . However, as the signature extraction of the weights is already performed, only the bias remains unknown. This implies that we only require one additional non-null critical point, corresponding to the targeted *always-on* neuron, to extract the bias.

Extraction of *input-off* neurons. Finally, the case of the *input-off* neurons has a very similar solution to the case of *always-on* neurons. The idea is also to use the signature’s search for a deeper neuron to extract the signature of the targeted neuron. However, for *input-off* neurons, their state changes with the input, and the activation function can no longer be approximated by the identity function for all inputs. Since we are targeting ReLU activations, we know that the activation function is equivalent to either the null or the identity function. Therefore, we use physical leakages to detect which inputs lead to the null or the identity function, and artificially transform the inputs’ equation ourselves. Then we solve the same system of equations as for the *always-on* neurons to extract the neuron’s signature. Similarly, to the case of *always-on* neurons, we extract the bias using a critical point directly for the targeted *input-off* neuron.

Our framework proposes a similar extraction method to what has been introduced in [6], [17], [28]. However, we use side-channel analysis both to make this methodology robust to special neurons, *e.g.*, *always-off/on* and *input-off*, and to improve the final part of our attack: extraction of the sign of the scaling factor. The following section describes how side-channel analysis can be used to infer the state of each neuron.

6. Stage 3: Sign Extraction

To achieve functional equivalence, the sign of the scaling factor associated with the extracted signature must be obtained in the third and final stage of our attack. We propose

a new methodology, free from any assumptions of access to the confidence scores or requiring supplementary queries, and only needing one hypothesis by layer. In Section 6.1, we introduce how we make use of the previously acquired physical traces to align all the neurons of the targeted layer. Further, in Section 6.2, we detail how we take advantage of this alignment to infer the scaling factor’s sign of all the neurons with only one hypothesis.

6.1. Layer alignment

The extraction of the sign of all the scaling factors begins once all of the neurons’ signature in the layer have been extracted. In this case, the attacker knows the induced hyperplane for each neuron and can associate an arbitrary state to each side of the hyperplanes, referred as V_A and V_B , creating 2^n hypotheses, with n the number of neurons in the layer. To reduce this number, we can align the neurons such as, for each neuron in the layer, the arbitrary state V_A corresponds to the same neuron’s real state (either active or inactive), while V_B corresponds to the other. The number of hypotheses is reduced to one corresponding to the identification of the sign related to V_A or V_B . This alignment is achieved using the collected physical traces during the signature’s search, *i.e.*, Stage 1. Indeed, after the first stage of our framework, for each neuron physical traces have been collected, labelled using the *k-means* algorithm and split into two groups corresponding to their arbitrary state V_A or V_B . Under the assumption that the processing of the ReLU function remains the same from one neuron to another within the same layer, it is possible to regroup the traces in two global groups by comparing the centroids from each cluster for each neuron.

6.2. Hypothesis selection

After the neurons of the targeted layer have been aligned, one assumption is made on the arbitrary state: we suppose that the state V_A (resp. V_B) corresponds to the active (resp. inactive) state. This hypothesis leads to two representations for the whole layer as represented in Figure 2. We test the validity of our hypothesis during the signature’s search of the first neuron in the next layer. Indeed, once enough critical points are extracted, a linear relation between these points and the output of the previous layer is expected (see Section 5). A wrong assumption on the state V_A and V_B will break this linear relation by having the ReLU function performing $A = ReLU(-x)$ instead of $A = ReLU(x)$. This produces incorrect activation values, *i.e.*, incorrect inputs for the next layer, and Equation 1 no longer holds for the next layer. Therefore, we use the residuals of the solution to the least-square problem (see Section 5) as a criterion for the validity of our sign hypothesis selection. The correct hypothesis is the one leading to the lowest residual error. Using this method, we are able to extract the sign of each neuron without requiring supplementary queries. Furthermore, since our method only requires one

hypothesis per layer, we argue that it remains practical even for complex and deep architectures.

7. Evaluation

In this section, we first describe the metrics that we use in Section 7.1 and the experimental setup in Section 7.2. Then, in Section 7.3, we detail the side-channel vulnerability we exploit in our methodology. Finally, we evaluate the practicability of our framework against architectures that have never been targeted before, such as a shortened MobileNetv1, in Section 7.4. We then compare our results with SOTA model extraction frameworks in Section 7.5.

7.1. Attack Evaluation Metrics

General metrics. Similarly to [6], we use the maximum of the absolute difference between the correct weights θ and their extracted counterparts $\hat{\theta}$, as well as the number of queries to evaluate the performances of our attack. We also use the mean error by layer between the true activation value A and the estimated one \hat{A} .

Metrics dedicated to classifiers. Most of the metrics described in [6] were designed for the specific task of regression. To evaluate our attack on other tasks, such as classification, we extend our set of metrics to include two additional ones used in [27] and [13]: *Fidelity* and *Accuracy Under Attack* (AuA). Fidelity corresponds to the average label agreement between the stolen and the original model on the *top-1* prediction for a subset of inputs. Accuracy under attack corresponds to the transfer rate of adversarial examples generated on the stolen model and tested on the targeted model. It is used as a proxy to represent the ability to craft more powerful attacks on the victim’s network using the stolen model.

7.2. Experimental setup

Dataset and Architecture. We target several embedded DNNs including a shortened MobileNetv1 [15]. This DNN shares the same basic block as the original MobileNetv1: depth-wise separable convolutions. To the best of our knowledge, these blocks have never been successfully extracted using cryptanalytic extraction. Full description of all the targeted architectures are provided in Appendix B. All classifiers were trained on the CIFAR-10 dataset¹ except for the architecture from [9]. This one, as well as all the DNNs performing regression, is trained on random datasets. Table 2 summarizes all the targeted DNNs in this work.

Hardware configuration. The models are embedded on an ARM Cortex-M7 microcontroller unit (MCU), on an STM32-F767ZI board via the X-Cube-AI framework, and uses 32-bit floating point data for all processed values (*i.e.*, input, weights, activations values, *etc.*). This board incorporates 2 Mbytes of flash memory and 512 Kbytes of

1. <https://www.cs.toronto.edu/~kriz/cifar.html>

TABLE 2: Targeted architectures in this work.

Type of architecture	Number of parameters	Number of targets	Targeted by
Small MLP	< 200,000	7	This work, [4], [6], [17], [28]
MLP	935,370	1	This work, [5]
Large MLP	1,721,802	1	This work
MobileNetv1-short	5,234	1	This work

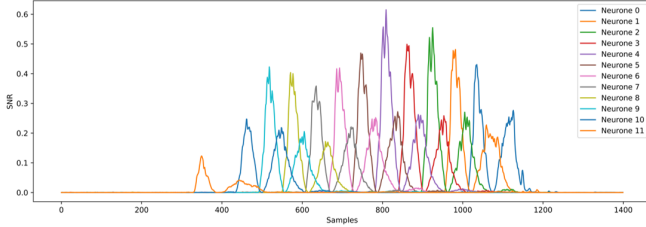


Figure 4: Signal-to-noise ratio associated with the processing of the ReLU function.

SRAM. X-Cube-AI is a framework developed by ST Microelectronics to facilitate the deployment of DNNs on STM’s embedded devices. However, it is important to notice that this framework is not intended to provide secure deployment for embedded DNNs.

Practical setup. A probe from Langer (EMV-Technik RF-U 2.5 is connected to an amplifier (ZLF-2000G+) to capture the EM physical leakages. To acquire the EM signal, a Lecroy oscilloscope (2.5 GHz WaveRunner 625Zi) is used in our setup.

7.3. Side-channel vulnerability

Our attack relies on the crucial step of side-channel extraction of each neuron’s state (see Section 4). To validate our methodology, we extract the state of each neuron from an embedded neural network, as described in Section 7.2.

Description of the weakness. The implementation of the ReLU function in X-Cube-AI is constant-time, as recommended in [24]. Algorithm 2 describes the main idea behind this implementation² for 32-bit values.

Algorithm 2 Pseudo-code of the constant-time ReLU function

Require: the vector of pre-activation values $\mathbf{Z} \in \mathbb{R}^n$ with n the number of neurons

Ensure: $\mathbf{A} = \max(0, \mathbf{Z}) \in (\mathbb{R}^+)^n$

for i in m **do**

$A_i \leftarrow (\sim (\mathbf{Z}_i \gg 31)) \& \mathbf{Z}_i$

end for

Given a signed pre-activation value $\mathbf{Z}_i \in \mathbb{R}$ encoded on 32 bits, the sign can be extracted ($\mathbf{Z}_i \gg 31$), and casted

2. For the purpose of this work, we decide to explain the weakness we exploit as an illustrated example, but understanding the vulnerability is not a prerequisite to perform our black-box extraction attack.

on 32 bits to serve as a mask (*i.e.*, $(\sim (\mathbf{Z}_i \gg 7)) \in \{0x00000000, 0xFFFFFFFF\}$). Hence, depending on the sign of \mathbf{Z}_i , application of the mask will either cause no modification or set the output to zero. A variant of this implementation can also be found in CMSIS-NN [20].

Although this implementation is constant-time, the usage of the mask leads to a physical leakage that can be captured by an attacker. In particular, depending on the state of the neuron, the mask is defined by $0x00000000$ or $0xFFFFFFFF$. The difference in the Hamming weight, *i.e.*, the number of bits equal to one between these two values, produces a data-dependency between the physical traces and the manipulated mask. As this dependence can be linked to the state of each neuron, it is therefore possible to create two groups, using clustering algorithms, on the physical traces, corresponding to the two active and inactive states.

Critical point search. To validate our method, we acquire a set of 10,000 EM traces using the setup defined in Section 7.2. Then, we conduct the strategy defined in Section 4 and measure the success rate of retrieving the neuron’s state for 12 neurons. To identify the dependence between the physical traces and the mask (*i.e.*, neuron’s state), we apply the Signal-to-Noise Ratio (SNR) [25] using the 10,000 traces. The obtained result is illustrated in Figure 4. As expected, a high dependence, *i.e.*, a high SNR, between the physical traces and the mask can be observed.

As mentioned in Section 4, we use the *k-mean* algorithm to generate the two neuron’s state clusters without any prior knowledge on the neuron’s states themselves. To measure the success rate, we suppose access to the true state of the neurons. As stated in Section 4, we only need a correct separation between the two clusters for our attack. Therefore, to compute our metric, we verify for each input, that if they correspond to the same state, they are in the same cluster.

With this method, we successfully extract the neuron’s state for the first time using a black-box side-channel attack on a protected implementation based on a SOTA countermeasure. We obtain an average success rate of 100% over all tested neurons and traces (see Appendix C for further details). This success rate and the number of physical traces needed to attain it is highly dependant of several experimental factors, such as the overall noise produced by the physical system, the type of the probe, *etc.*. Given the high number of queries needed to fully extract a DNN, it is critical to maintain a high success rate while keeping the number of needed traces as low as possible. To consider the worst case scenario for the victim (*i.e.*, best case scenario for the attacker), we suppose an ideal set-up allowing the attacker to perform side-channel attacks with a 100% success rate in only one EM trace. All of the followings results are given under this hypothesis. This allows to set an empirical boundary for the performance of our framework and will ease future comparison with other frameworks. Finally, this assumption is not restrictive since, if one attack requires N EM traces to achieve a 100% success rate, the global complexity of the attack is given by the worse-case scenario complexity multiplied by N .

Sign extraction. As previously mentioned, the side-

TABLE 3: Results of our extraction framework on a large MLP and a truncated MobileNetv1.

Model	Parameters	Queries	Accuracy Original	Accuracy Stolen	Fidelity
MLP	1, 721, 802	$2^{26.0}$	54.3%	54.0%	93.2%
MobileNetv1	5, 234	$2^{18.8}$	59.7%	59.2%	88.4%

channel vulnerability is not only used to find the critical points during Stage 1 of our attack, but also to determine the sign of the neurons once all of the neuron’s signatures in the layer have been found. For this step, we exploit the methodology introduced in Section 6 to reduce the number of clusters to two for the whole layer. Based on this method, we achieve a success rate of 99.7% for retrieving the sign of the scaling factor related to each neuron.

Our results are not dependent from the algorithm, and could be improved using more complex unsupervised methods, but the overall principle will remain the same. To the best of our knowledge, this is the first successful extraction of the neuron’s sign in a hard-label setting not requiring extra queries.

7.4. Evaluation of the framework against large architectures

We first consider the extraction of a large MLP. While this architecture could be targeted by other SOTA frameworks, as it is only composed of fully connected layers, we use it to validate our framework against a very large architecture with almost twice the number of weights of the largest DNN targeted by previous cryptanalytic extraction frameworks [5]. Then, to test the true advantage of our framework, *i.e.*, extraction of non-fully connected layers, we target a truncated MobileNetv1. While its number of parameters is quite low (5, 234) compared to the extracted MLP, this is, to the best of our knowledge, the deepest DNN extracted using cryptanalytic methods. Furthermore, this model includes depth-wise separable convolutions followed by batch-normalization, which has never been targeted before. It also includes an average pooling layer, which makes it impossible to use the methods introduced in [6], [17], [28] or in [5].

We summarize our results in Table 3. We successfully extract a non-fully connected DNN with high fidelity for the first time using our framework and achieve an accuracy of 59.2% (resp. 54.9%) with a fidelity of 88.4% (resp. 94.1%) on the complete test-set of CIFAR-10 (10, 000 images) for the MobileNetv1 (resp. large MLP).

Comparison with learning-based extraction. We provide a comparison with learning-based extractions using a baseline model trained on the queries sent during the critical points’ search, with the responses of the targeted model as labels. To be in the most ideal settings for the baseline model, the same hyper-parameters (architecture, learning rate, *etc.*) as for the training of the targeted model are used, as well as a balanced training dataset to obtain similar number of samples per class. While the baseline

method seems to offer decent performance during training (around 56% of accuracy on the random dataset), it does not generalize well and the performance drops significantly on the true dataset, achieving only 19.6% with a fidelity of 21.1% for the MobileNetv1. The use of random queries, which are not from the same distribution as the true training set from CIFAR-10, and the hard-label setting drastically reduce the performance of the baseline model. This illustrates that learning-based extraction methods are limited by data distribution which makes cryptanalytic attacks more appropriate for offering better performance in the threat model considered in this paper.

Impact of the depth. Using binary search, it is only possible to place a critical point within a certain distance of the hyperplane (see Section 4.2). This creates a small error on the estimation of the weights as we can see in Table 4. For the first layer L0, the maximum error on the estimated weights is $2^{-18.9}$. As we exploit side-channel attacks to localize critical points, we are primarily limited by the data type processed by the DNN. For comparison, we include in Table 4 the results of our extraction framework on the same target but using 64-bit data. In this scenario, the maximum error for the same layer L0 is reduced to $2^{-46.6}$ which corresponds to a division of the maximum error by a factor over 500 compared to SOTA best results in a similar scenario [6]. This proves the value of side-channel attacks for the critical point search. However, from Table 4, we can also see that the performance of the extraction attack diminishes with the depth of the layer. This is due to the propagation of this small error on the weights in the current layer, which will be accumulated in the activation values, *i.e.*, the inputs for the system of equations of the next layer. This observation is not intrinsic to our method, as all SOTA cryptanalytic extraction frameworks are impacted by it. It is then crucial to minimize the error at each layer to have the smallest impact on the next layer, which underlines the benefits of our extraction. Alternatively, several methods could be implemented to reduce this propagation, such as re-computation of the weights to achieve higher precision, but we leave them for future works.

Impact of the last layer. As we can see in Table 4, the mean difference between the activation values from the victim’s model and the one from our stolen model remains below 0.01 even after the extractions of the first eleventh layers. Therefore, we suppose that the loss in the fidelity between the stolen and the targeted model comes mainly from the extraction of the last layer. To validate our assumption, we measure the fidelity and the accuracy obtained from an hybrid MobileNetv1 model composed of our first eleventh layers and a final layer with the real weights from the targeted model. Such model achieves a *top-1* accuracy of 59.77% on CIFAR10, and a fidelity of 99.56%. This confirms that the main loss in the fidelity is coming from the extraction of the last layer. The absence of ReLU function in this layer implies that no critical points can be extracted for this layer. Therefore, we use the collected input-output pairs during Stage 1, and supervised learning to extract the last layer’s weights. While this method is not optimal, the high

TABLE 4: Layer wise statistics of the extraction of the shortened MobileNetv1.

Datatype	Metrics	L0	L1	L2	L3	L4	L5	L6	L7	L8	L9	L10	L11
32 bits	$\max \theta - \hat{\theta} ^L$	$2^{-18.9}$	$2^{-17.6}$	$2^{-7.9}$	$2^{-18.2}$	$2^{-7.6}$	$2^{-13.9}$	$2^{-9.9}$	$2^{-11.4}$	$2^{-8.8}$	$2^{-6.7}$	$2^{-4.3}$	$2^{3.7}$
	$\text{mean} A - \hat{A} ^L$	$2^{-21.6}$	$2^{-20.9}$	$2^{-11.2}$	$2^{-11.2}$	$2^{-15.2}$	$2^{-15.5}$	$2^{-13.0}$	$2^{-12.7}$	$2^{-12.3}$	$2^{-11.7}$	$2^{-8.6}$	$2^{4.1}$
64 bits	$\max \theta - \hat{\theta} ^L$	$2^{-46.6}$	$2^{-43.8}$	$2^{-37.4}$	$2^{-34.2}$	$2^{-29.0}$	$2^{-27.1}$	$2^{-26.0}$	$2^{-26.8}$	$2^{-23.1}$	$2^{-22.7}$	$2^{-15.3}$	$2^{3.8}$
	$\text{mean} A - \hat{A} ^L$	$2^{-49.0}$	$2^{-48.0}$	$2^{-40.3}$	$2^{-34.1}$	$2^{-32.2}$	$2^{-31.6}$	$2^{-28.3}$	$2^{-28.0}$	$2^{-26.4}$	$2^{-25.6}$	$2^{-20.5}$	$2^{4.4}$

precision achieved in the extraction of the previous layer outweighs the induced loss which allows our extraction to offer high fidelity. As this limitation impacts all methods based on extraction of the critical points in hard-label settings, we leave improvement of the extraction of the last layer for future works.

Generation of adversarial examples. Finally, we use AuA to quantify how much our extraction attack can improve other attacks, and more specifically adversarial attacks in the settings described in [13], [27]. We obtain a transfer rate of 95.78% (resp. 96.67%) for the MobileNetv1 (resp. large MLP) architecture, achieving close to white-box performance. Our methodology outperforms powerful substitute model trained with access to the true data [27] and prove its ability to mount more powerful attacks on the targeted model.

We provide supplementary results on the MLP and MobileNetv1 architectures in Appendix D as well as results on the impact of the special cases of neurons in Appendix E. We then compare our framework against results from SOTA methodologies in the next section.

7.5. Evaluation of the framework against SOTA techniques

Since cryptanalytic extraction of DNNs running 32-bit data has not been previously performed, only results with 64-bit data are provided in this section to have a more representative comparison. The results we present are obtained via simulating the side-channel vulnerability on software implementation with DNNs using 64-bit data. These results are synthesized in Table 5 To get the most complete work, we also conduct our extraction attack on 32-bit and report the result in Appendix F.

TABLE 5: Simulation of our framework on DNN with 64-bit data.

Architecture	Parameters	Approach	Queries	$\max \Delta_{\theta} ^{L-1}$
784-128-1 (R)	100, 480	This work	$2^{22.6}$	$2^{-40.8}$
		[6]	$2^{21.5}$	$2^{-24.7}$
40-20-10-10-1 (R)	1, 110	This work	$2^{16.8}$	$2^{-42.0}$
		[6]	$2^{17.8}$	$2^{-27.1}$
3072-256-256-256-64-10 (C)	935, 370	This work	$2^{26.2}$	$2^{-35.0}$
		[5]	-	$2^{-18.0}$

(C) Classification task with hard-label settings, (R) Regression task with access to the prediction score.

The 784-128-1 DNN corresponds to a model with an input size of 784, an hidden layer of 128 neurons and 1 output. Except for the work included in [5], all the other frameworks consider access to the confidence scores and exploit

the estimated gradient. To have a fair comparison for the cryptanalytic frameworks, we excluded the final layer of our results in Table 5. Full results are included in Appendix F.

Side-channel attacks offer a higher precision on the localisation of the critical points than previous SOTA methods. In consequence, our framework improves accuracy on the extraction of the DNNs targeted in [6] by reducing the maximum error on the extracted weights from $2^{-24.7}$ (resp. $2^{27.1}$ to $2^{-40.8}$ (resp. $2^{-42.0}$) corresponding to a division by more than 30,000 of this error. For the comparison with the framework presented in [5], as they operate under the assumption that the previous layer is perfectly extracted, we also use this assumption, and reduce the maximum error on the estimated weights from $2^{-18.0}$ to $2^{-35.0}$. Supplementary results are given in Appendix F. In [5], the dual point’s search is also performed via binary search. Given the limited level of information provided by the hard-label setting, this increases the number of necessary queries, making our attack more effective. However, the precise number of queries was not given so far, so the comparison is rather hard. Also, when access to the gradient is authorized, the gain we obtain from the direct association between the critical points and their neuron compensates the loss induced by the naive binary search (*i.e.*, without optimization using the gradient [6]).

We give supplementary results against SOTA frameworks as well as a more detailed analysis of our results on the 3072-256-256-256-64-10 architecture, in Appendix F. In conclusion, we can see that our method is able to extract with high fidelity all of the architectures targeted by other SOTA frameworks, and even outperforms them on the precision obtained on the extracted weights.

7.6. Related works

In this work, we only compare ourselves to the results obtained with the cryptanalytic extraction frameworks presented in [4], [6], [17], [28] and in [5]. We argue that the difference in threat model and attacker’s goals for the extraction is too high to make a relevant comparison with other type of frameworks.

Side-channel extraction of DNNs using CEMA, as presented in [2], [14], naturally offers guaranties in term of fidelity, but suffers from several limitations, as mentioned in [19]. In particular, the extraction of the bias or the impact of the ReLU function on the attack are still open problems, which limit the practicability of these attacks. Similarly to our framework, Zhang *et al.* also use side-channel attacks to infer the state of the neurons and the confidence scores to extract each layers [33]. However, due to their access to

the Intel Running Average Power Limit (RAPL) and to an open device to perform their profiled side-channel attack, they consider a very different threat model to ours.

Model extraction frameworks using fault injection have been recently proposed [13], [27]. Both have proven successful against a variety of architectures using the Rowhammer attack and laser fault injection, respectively. However, the threat model considered in [13] is very close to the one in [33] *i.e.*, access to the logits and to an open device, making the comparison with our work biased. Finally, the work in [27] requires that the attacker and the victim share the same device, *i.e.* MLaaS context, and is only applicable on DRAM platforms which corresponds to a very different threat model from our own.

8. Conclusion

In this study, we propose a new framework combining side-channel attacks and cryptanalytic extraction to target fully and non-fully connected DNNs. We successfully extract complex architectures with high fidelity, emphasizing on the serious threat that poses side-channel attacks on the intellectual property (IP) of DNNs. In particular, the lack of protection on the activation functions from a side-channel context has allowed the extraction of the weights of a shortened MobileNetV1. Such vulnerability could lead to different exploitation in the future and could easily be adapted to adversarial attacks for example. Therefore, it is critical to integrate both software and hardware countermeasures against model extraction to protect the IP of DNNs.

References

- [1] R.o. duda, p.e. hart, and d.g. stork, pattern classification, new york: John wiley & sons, 2001, pp. xx + 654, isbn: 0-471-05669-3. *Journal of Classification*, 24:305–307, 2007.
- [2] 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.
- [3] E. Cagli, C. Canovas, and E. Prouff. Convolutional neural networks with data augmentation against jitter-based countermeasures - profiling attacks without pre-processing. In *Workshop on Cryptographic Hardware and Embedded Systems*, 2017.
- [4] I. A. Canales-Martínez, J. Chávez-Saab, A. Hambitzer, F. Rodríguez-Henríquez, N. Satpute, and A. Shamir. Polynomial time cryptanalytic extraction of neural network models. In M. Joye and G. Leander, editors, *Advances in Cryptology – EUROCRYPT 2024*, pages 3–33, Cham, 2024. Springer Nature Switzerland.
- [5] N. Carlini, J. Chávez-Saab, A. Hambitzer, F. Rodríguez-Henríquez, and A. Shamir. Polynomial time cryptanalytic extraction of deep neural networks in the hard-label setting. 2024.
- [6] N. Carlini, M. Jagielski, and I. Mironov. Cryptanalytic extraction of neural network models. In *Annual International Cryptology Conference*, 2020.
- [7] N. Carlini and D. A. Wagner. Towards evaluating the robustness of neural networks. *2017 IEEE Symposium on Security and Privacy (SP)*, pages 39–57, 2016.
- [8] P.-Y. Chen, H. Zhang, Y. Sharma, J. Yi, and C.-J. Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, 2017.
- [9] Y. Chen, X. Dong, J. Guo, Y. Shen, A. Wang, and X. Wang. Hard-label cryptanalytic extraction of neural network models. *IACR Cryptol. ePrint Arch.*, 2024:1403, 2024.
- [10] M. Fredrikson, S. Jha, and T. Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015.
- [11] Y. Gao, H. Qiu, Z. Zhang, B. Wang, H. Ma, A. Abuadba, M. Xue, A. Fu, and S. Nepal. Deeptheft: Stealing dnn model architectures through power side channel. In *2024 IEEE Symposium on Security and Privacy (SP)*, pages 3311–3326, 2024.
- [12] B. Hanin and D. Rolnick. Deep relu networks have surprisingly few activation patterns. In *Neural Information Processing Systems*, 2019.
- [13] K. Hector, M. Dumont, P.-A. Moellic, and J.-M. Dutertre. Fault injection and safe-error attack for extraction of embedded neural network models. In *ESORICS 2023 International Workshops*, volume 14399 of *Computer Security. ESORICS 2023 International Workshops CPS4CIP, ADIoT, SecAssure, WASP, TAURIN, PriST-AI, and SECAI, The Hague, The Netherlands, September 25–29, 2023, Revised Selected Papers, Part II*, pages 644–664, La Hague, Netherlands, Sept. 2023. Delft University of Technology.
- [14] P. Horvath, L. Chmielewski, L. Weissbart, L. Batina, and Y. Yarom. Barracuda: Bringing electromagnetic side channel into play to steal the weights of neural networks from nvidia gpus. *ArXiv*, abs/2312.07783, 2023.
- [15] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *ArXiv*, abs/1704.04861, 2017.
- [16] S. H. Huang, N. Papernot, I. J. Goodfellow, Y. Duan, and P. Abbeel. Adversarial attacks on neural network policies. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. OpenReview.net, 2017.

- [17] M. Jagielski, N. Carlini, D. Berthelot, A. Kurakin, and N. Papernot. High accuracy and high fidelity extraction of neural networks. In *USENIX Security Symposium*, 2019.
- [18] R. Joud, P. Moëllic, S. Pontié, and J. Rigaud. Like an open book? read neural network architecture with simple power analysis on 32-bit microcontrollers. In S. Bhasin and T. Roche, editors, *Smart Card Research and Advanced Applications - 22nd International Conference, CARDIS 2023, Amsterdam, The Netherlands, November 14-16, 2023, Revised Selected Papers*, volume 14530 of *Lecture Notes in Computer Science*, pages 256–276. Springer, 2023.
- [19] R. Joud, P.-A. Moëllic, S. Pontié, and J.-B. Rigaud. A practical introduction to side-channel extraction of deep neural network parameters. In *Smart Card Research and Advanced Application Conference*, 2022.
- [20] L. Lai, N. Suda, and V. Chandra. Cmsis-nn: Efficient neural network kernels for arm cortex-m cpus. *ArXiv*, abs/1801.06601, 2018.
- [21] J. Ma. A higher-level Neural Network library on Microcontrollers (NNoM), Oct. 2020.
- [22] H. Maghrebi, T. Portigliatti, and E. Prouff. Breaking cryptographic implementations using deep learning techniques. In *Security, Privacy, and Applied Cryptography Engineering: 6th International Conference, SPACE 2016, Hyderabad, India, December 14-18, 2016, Proceedings 6*, pages 3–26. Springer, 2016.
- [23] H. T. Maia, C. Xiao, D. Li, E. Grinspun, and C. Zheng. Can one hear the shape of a neural network?: Snooping the gpu via magnetic side channel. In *USENIX Security Symposium*, 2021.
- [24] S. Maji, U. Banerjee, and A. P. Chandrakasan. Leaky nets: Recovering embedded neural network models and inputs through simple power and timing side-channels—attacks and defenses. *IEEE Internet of Things Journal*, 8:12079–12092, 2021.
- [25] S. Mangard, E. Oswald, and T. Popp. *Power analysis attacks - revealing the secrets of smart cards*. Springer, 2007.
- [26] S. Milli, L. Schmidt, A. D. Dragan, and M. Hardt. Model reconstruction from model explanations. *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 2018.
- [27] A. S. Rakin, M. H. I. Chowdhury, F. Yao, and D. Fan. Deepsteal: Advanced model extractions leveraging efficient weight stealing in memories. *2022 IEEE Symposium on Security and Privacy (SP)*, pages 1157–1174, 2021.
- [28] D. Rolnick and K. P. Kording. Reverse-engineering deep relu networks. In *International Conference on Machine Learning*, 2019.
- [29] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- [30] J. Uesato, B. O’donoghue, P. Kohli, and A. Oord. Adversarial risk and the dangers of evaluating against weak attacks. In *International conference on machine learning*, pages 5025–5034. PMLR, 2018.
- [31] Y. Umuroglu, N. J. Fraser, G. Gambardella, M. Blott, P. Leong, M. Jahre, and K. Vissers. Finn: A framework for fast, scalable binarized neural network inference. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays, FPGA ’17*, page 65–74, New York, NY, USA, 2017. Association for Computing Machinery.
- [32] H. Yu, H. Ma, K. Yang, Y. Zhao, and Y. Jin. Deepem: Deep neural networks model recovery through em side-channel information leakage. In *2020 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*, pages 209–218, 2020.
- [33] X. Zhang, A. A. Ding, and Y. Fei. Deep-learning model extraction through software-based power side-channel. *2023 IEEE/ACM International Conference on Computer Aided Design (ICCAD)*, pages 1–9, 2023.

Appendix

Appendix A.

Approximation of the *always-on* and *input-off* neurons

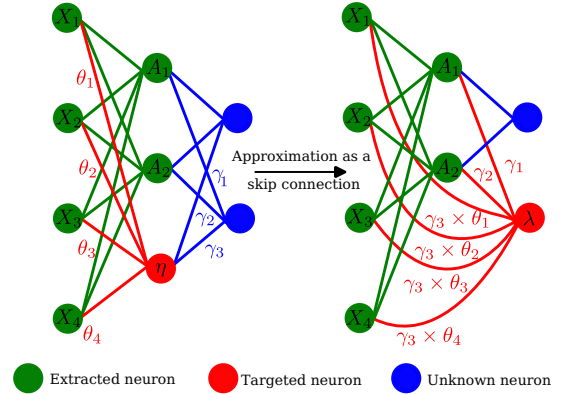


Figure 5: Approximation of the special case neuron η as a skip connection to the neuron λ .

In Figure 5 we provide a schematic representation of a targeted *always-on* neuron being extracted in the next layer. The neuron η is approximated as a weighted skip connection to the next layer. It is then possible to extract its signature by finding the signature of a neuron in the following layer. In the case represented in this figure, we use the neuron λ in the next layer to achieve this. As mentioned in subsection 5.2, using this method, the signature of λ includes also the signature of η as described in Equation 5. This trick takes advantages of the ReLU function to approximate it to the identity function for positive values.

$$\sum_{i=1}^2 \gamma_i A_i + \gamma_3 \times \sum_{j=1}^4 \theta_j X_j + \gamma_3 \beta_\eta + \beta_\lambda = 0, \mathbf{X} \in \mathbb{R}^4$$

Appendix B.

Description of the targeted architectures

An overall description of the targeted architecture based on MobileNetv1 is given in Table 6 and Table 7. Then in Table 8, we detail the architecture of our large MLP.

Appendix C.

Supplementary results on the side-channel attack

In our experiment we use the *k-means* algorithm to infer the state of the neurons in our search for the critical points. We achieve an average success rate of 85.9%, using only one side-channel trace for each acquisition. While this result validates our methodology, the high number of queries

TABLE 6: Description of the architecture of the shortened MobileNetv1.

Layers	Output size	Parameters
Conv2D 3×3 , padding, stride=2	$1 \times 12 \times 16 \times 16$	324
BatchNorm2D	$1 \times 12 \times 16 \times 16$	24
ReLU	$1 \times 12 \times 16 \times 16$	0
Depth-wise separable convolution	$1 \times 8 \times 16 \times 16$	244
Depth-wise separable convolution	$1 \times 16 \times 8 \times 8$	248
Depth-wise separable convolution	$1 \times 16 \times 8 \times 8$	464
Depth-wise separable convolution	$1 \times 32 \times 4 \times 4$	752
Depth-wise separable convolution	$1 \times 64 \times 4 \times 4$	2528
Adaptive Average Pooling	$1 \times 64 \times 1 \times 1$	0
Dense	1×10	650

TABLE 7: Description of the depth-wise separable convolution.

Layers	Kernel size	Padding	Grouped convolution	Bias
Conv2D BatchNorm2D ReLU	3×3	1	✓	✗
Conv2D BatchNorm2D ReLU	$N_c \times 1 \times 1$	1	✗	✗

N_c corresponds to the number of output channel of the previous layer.

needed to perform the attack (see Table 3) imposes a higher success rate (close to 100 %) to extract complex DNNs. A basic method to increase the success rate, would be to repeat the acquisition of the EM trace and use a majority vote to infer the state.

In Figure 6, we plot the success rate obtained with a majority vote versus the number of times the side-channel acquisition is repeated for different initial success rate (*i.e.*, success rate in one trace). With a 85.9% initial success rate, we achieve a true success rate over 99% with 7 traces and over 99.9% with 13 traces. In practice, it is possible to increase the number of side-channel traces to come close to

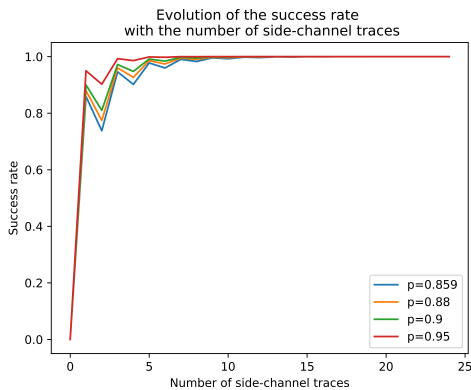


Figure 6: Success rate vs number of side-channel traces.

TABLE 8: Description of the architecture of the large MLP.

Layers	Output size	Parameters
Flatten	1×3072	0
Dense + ReLU	1×512	1,573,376
Dense + ReLU	1×256	131,328
Dense + ReLU	1×64	16,448
Dense + ReLU	1×10	650

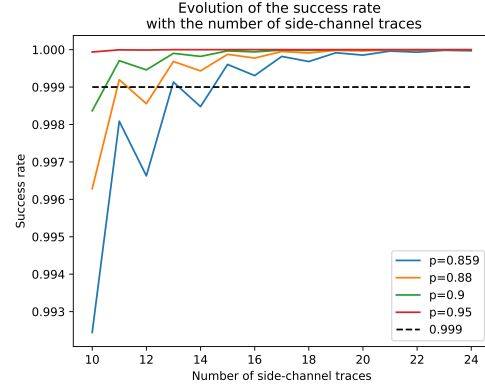


Figure 7: Success rate vs number of side-channel traces (zoom of Figure 6).

a 100% success rate with an arbitrary precision with a rate of convergence bounded by the Chernoff Bound. Furthermore, the global success rate is highly impacted by the success rate in one trace, which is highly dependent on the implementation and the experimental set-up. This success rate could greatly be improved in an ideal scenario where the chip would be unpackaged, which was not the case in our set-up. Finally, the choice of the algorithm for the classification of the side-channel trace has an important impact on the final result. For all of the previously mentioned reasons, we choose to focus our work on the methodological aspects and give our results considering an ideal attack. We leave the improvement of the side-channel attacks for future work.

Appendix D. Supplementary results for the MLP and the MobileNetv1 architecture.

In the MobileNetv1 model, the 10 successive layers after the first convolution layer (denoted L1-L10 in Table 4) correspond to the depth-wise separable convolutions. These blocks are composed of two convolutions, one depth-wise with a kernel of $1 \times 3 \times 3$ weights and, a point-wise convolution composed of $c \times 1 \times 1$ weights, with c the number of channels. The alternation between these two types of convolution is clearly visible in Table 4, the point-wise convolution corresponds to the even layers and the depth-wise to the odd layers. The extraction for the depth-wise convolution is always easier than for the point-wise. This is likely due to the number of parameters being lower for the depth-wise. However, it is not the only factor, as for layer L4, the number of channels is only of 8. And, as such, the

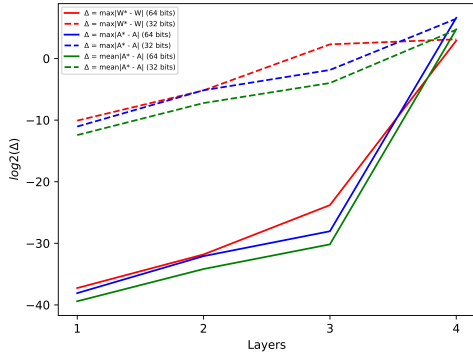


Figure 8: Results on the MLP architecture presented in Table 8.

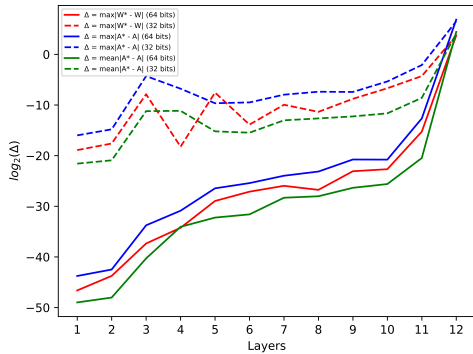


Figure 9: Results on the shortened MobileNetV1 architecture.

number of parameters is lower than for the corresponding depth-wise. But the precision of the extraction remains lower for the point-wise compared to the depth-wise. We make the hypothesis that the higher error comes from the different operations induced in the point-wise convolution in comparison with depth-wise convolution. In particular, for the point-wise convolution, the error could be higher because its operations mix the channels, whereas the depth-wise convolution performs its operations on a single channel. This behaviour is not present in the results of the extraction of the MLP.

Appendix E. Impact of the special cases neurons

To visualize the impact of the special cases of neurons, we present their share in the total number of requests and neurons for the 3072-256-256-256-64-10 MLP in Figure 10.

In Figure 10, we display the percentage of requests dedicated to the special cases neurons by layer as well as the percentage of these neurons by layer. As expected, the number of special neurons augments with the depth of

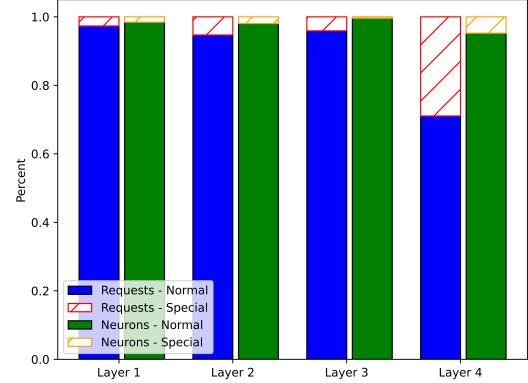


Figure 10: Comparison of the percentage of request dedicated to normal versus special cases neurons.

TABLE 9: Average number of requests by layer for normal and special cases neurons.

Neuron type	Layer 1	Layer 2	Layer 3	Layer 4
Special	32, 919.8	53, 158.8	224, 354.0	416, 799.0
Normal	19, 051.2	18, 885.1	20, 634.9	51, 136.7

the layer. This observation is confirmed when looking at the average number of requests to extract one special case neuron and one normal neuron for each layer (see Table 9). Among those special cases neurons, most of them can be defined as *input-off*. Due to the nullification of some portions of the input space by ReLU functions, obtaining a full rank matrix for the system of equations is difficult which leads to an increase of the number of requests. This issue is more present in deeper layer due to nesting of the activation functions which restricts the input space accessible to an attacker in deeper layer. One global consequence of this nesting is the increase in the average number of queries necessitated for one neuron, independently of its type (normal or special case) extraction. These results highlight the impact of the special cases neurons on the number of requests to extract a DNN even for a simple architecture (*i.e.* MLP).

Appendix F. Supplementary results against SOTA frameworks

To further validate our contribution, we compare the results obtained with our method against the ones obtained with other SOTA frameworks. However, since most of these frameworks (except for the one included in [5], [9]) target DNNs trained for regression tasks, the unrestricted access to the output value of the DNN allows the attacker to compute the gradient. Furthermore, all frameworks consider DNN using 64-bit data, which allows for an accurate estimation of the gradient. We decided to compare our framework in that context, the overall methodology does not change, but access to the continuous output value improves our extraction of the last layer (see Section 7.4). Our results are presented in Table 10.

TABLE 10: Simulation of our framework on DNN with 64-bit data (supplementary results).

Architecture	Parameters	Approach	Queries	$\max \Delta_\theta ^L$
784-32-1	25,120	This work	$2^{20.6}$	$2^{-43.5}$
		[6]	$2^{19.2}$	$2^{-30.2}$
		[17]	$2^{18.2}$	$2^{-1.7}$
10-10-10-1	210	This work	$2^{15.6}$	$2^{-46.2}$
		[6]	$2^{16.0}$	$2^{-36.0}$
		[28]	$2^{22.0}$	$2^{-12.0}$
10-20-20-1	620	This work	$2^{15.6}$	$2^{-46.5}$
		[6]	$2^{17.1}$	$2^{-37.0}$
80-40-20-1	4,020	This work	$2^{18.3}$	$2^{-44.2}$
		[6]	$2^{18.5}$	$2^{-39.7}$

We outperform all the other cryptanalytical frameworks in the maximum distance between the estimated weights and the real ones. Even in an ideal side-channel scenario (see Section 7.3 and Appendix C), our attack requires a higher number of queries for short and wide network than other SOTA attacks. However, if we increase the depth of the targeted network, we can see that the performance of our attack improves compared to other frameworks. This could result from an optimization described in [6], where the gradient is used to localize directly the critical points without binary search, in the case where only one neuron changes sign. In this configuration, their search is greatly improved compared to ours. The counterpart is that, since the gradient does not provide an association between the critical points and the neurons, the side-channel does. In order to be sure to have enough critical points for each neuron, gradient-based methods need to extract a higher number of critical points than side-channel based methods do. This trade-off could explain our performance and why the results obtained with our framework are better in the case of deeper DNNs. To facilitate future comparisons, we also include results on all these architectures but using 32-bit data in Table 11.

TABLE 11: Results of our framework against targeted architecture by SOTA frameworks using 32-bit data

Architecture	Parameters	Queries	$\max \Delta_\theta ^L$
784-32-1	25,120	$2^{19.8}$	$2^{-17.7}$
784-128-1	100,480	$2^{21.7}$	$2^{-17.4}$
10-10-10-1	210	$2^{13.0}$	$2^{-18.2}$
10-20-20-1	620	$2^{14.5}$	$2^{-17.8}$
40-20-10-10-1	1,110	$2^{16.4}$	$2^{-12.1}$
80-40-20-1	4,020	$2^{19.1}$	$2^{-14.8}$

We also include comparison with recent works by Chen *et al.* [9] and Carlini *et al.* [5]. In both studies, the authors demonstrate that the attacker could target DNNs in a hard-label settings. We thus present the full results obtained with our framework against the architectures targeted in these papers in Table 12. We include only the deepest architecture from [9], and we use the metrics associated with classifier similarly to Table 3. In Table 3, the metric $\max|\Delta_\theta|^L$ computes the maximum error between the estimated weights and the real ones. However, as mentioned in Section 7.4, the extraction of the last layer cannot be performed using our framework, and we therefore use a learning based ap-

proach to extract the weights of this layer. In Figure 11 and Figure 12, we consider this metrics over the whole network. We can clearly see that, for the previous layers, we propose a very low error on the weights or the activation values, and even outperform SOTA frameworks on the error between the estimated and the true weights. Once again, this highlights the need for an alternative method for the extraction of the last layer.

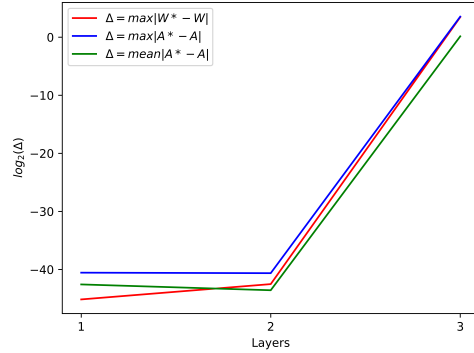


Figure 11: Results on the architecture presented in [9].

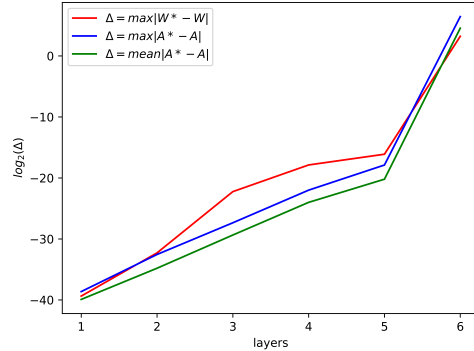


Figure 12: Success metrics over the whole network with our framework.

We also outperform methods presented in [5], [9] in terms of number of queries to the targeted DNN. In the hard-label settings this was expected, since we use side-channel to guide our binary search towards the critical point, whereas the works presented in [5], [9] use the boundary decision, which leaks fewer information on the state of the neurons. We achieve high fidelity between the stolen DNN and the targeted one against both architecture.

In [5], a perfect extraction of the previous layers is assumed before extraction of the current layer. This is a very strong assumption as it completely removes the propagation of any error. For comparison, we perform the extraction of the 3072-256-256-256-64-10 MLP under this assumption and report the results in Figure 13. We exclude the last layer in these results, as it not extracted using cryptanalyt-

TABLE 12: Comparison with state-of-the-art frameworks against classifiers.

Model	Parameters	Approach	Queries	Fidelity	$\max \Delta_\theta ^{L-1}$	$\max \Delta_\theta ^L$
1024-2-2-1	2,059	[9]	$2^{22.49}$	100%	-	$2^{-20.38}$
		This work	$2^{17.12}$	99.5%	$2^{-42.5}$	$2^{3.45}$
3072-256-256-256-64-10	935,370	[5]	-	-	-	$2^{-18.0}$
		This work	$2^{26.20}$	96.5%	$2^{-16.1}$	$2^{3.22}$

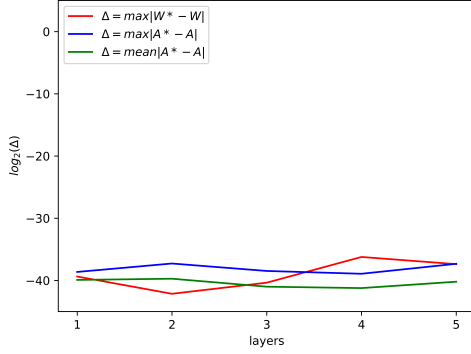


Figure 13: Extraction of the architectures presented in [5] without the propagation of errors and excluding the last layer.

ical methods. Under this hypothesis, we can see that the performance is not impacted by the depth of the DNN, and the extraction proposed by of our framework offers similar performance at each layer. This clearly underlines the importance of minimizing the propagation error during the extraction of DNNs.