

The Need for Speed: A Fast Guessing Entropy Calculation for Deep Learning-based SCA

Guilherme Perin, Lichao Wu, and Stjepan Picek *Senior Member, IEEE*

Abstract—The adoption of deep neural networks for profiling side-channel attacks (SCA) opened new perspectives for leakage detection. Recent publications showed that cryptographic implementations featuring different countermeasures could be broken without feature selection or trace preprocessing. This success comes with a high price: extensive hyperparameter search to find optimal deep learning models. As deep learning models usually suffer from overfitting due to their high fitting capacity, it is crucial to avoid over-training regimes, which require a correct number of epochs. For that, *early stopping* is employed as an efficient regularization method that requires a consistent validation metric. Although guessing entropy is a highly informative metric for profiling SCA, it is time-consuming, especially if computed for all epochs during training and the number of validation traces is significantly large. This paper shows that guessing entropy can be efficiently computed during training by reducing the number of validation traces without affecting the efficiency of early stopping decisions. Our solution significantly speeds up the process, impacting hyperparameter search and overall profiling attack performances. Our fast guessing entropy calculation is up to $16\times$ faster, resulting in more hyperparameter tuning experiments and allowing security evaluators to find more efficient deep learning model.

Index Terms—Side-channel Analysis, Deep learning, Guessing entropy, Validation phase, Fast Guessing Entropy.

1 INTRODUCTION

Side-channel attacks (SCA) explore the unintentional leakages (power consumption, time, and electromagnetic emissions) from electronic devices running secret-sensitive operations such as embedded cryptographic algorithms. Profiling SCA [1], [2], one of the most popular attack methods, is widely considered by developers and manufacturers when assessing worst-case security with strongest adversary assumptions. This type of attack assumes an adversary has a clone (open) device to build the strongest possible probabilistic model from collected side-channel measurements. Thus, the adversary applies the model to the victim's devices to recover the secret. If the profiling model is correct and can learn existing side-channel leakages, a profiling attack phase usually requires fewer side-channel measurements in comparison to non-profiling attacks [3]–[5].

Template attacks are the most classic form of profiling SCA [1]. Template attacks theoretically represent the strongest profiling model because of the typical underlying statistical distribution of side-channel leakages following multivariate Gaussian (or normal) distributions. Machine learning methods have also been considered for profiling attacks [6], [7], while their statistical parameters are learned from side-channel measurements rather than directly computed. Both Gaussian templates and machine learning models require feature selection. In the case of protected cryptographic implementations, the inability to make efficient feature selection (by selecting leakage samples with the highest Signal-to-Noise Ratios) may become a limiting factor to building optimal profiling models. Indeed, an effective

feature selection requires a strong correlation between the leakages and processed intermediate data. For instance, in the presence of masking countermeasures, evaluators select the best features (or points of interest) by knowing the secret random masks. Then, one can deploy worst-case security evaluations to emulate the adversaries having access to source code and secret shares during profiling. Additionally, template and machine learning-based models are susceptible to desynchronization effects in side-channel measurements, thus bringing additional challenges.

In recent years, the adoption of deep neural networks (DNNs) for profiling SCA provided competitive (and, in some cases, superior) results compared to template attacks and classical machine learning-based methods, especially against AES implementations [8], [9]. Without feature selection, which implies considering a weaker adversary, deep learning-based SCA can break cryptographic implementations protected with different countermeasures, such as Boolean masking and timing desynchronization [10]. Their high complexity follows the high learning capacity of DNNs; the expensive hyperparameter tuning becomes a limitation to fully explore the full potential of deep learning to find vulnerabilities in software and hardware implementations.

To make the hyperparameter tuning process more efficient, one tries to define appropriate hyperparameters ranges, which directly reflect the number of trainable parameters. Indeed, smaller DNNs may limit the learning capacity of a model, underfitting the profiling side-channel traces and providing poor attack performance. On the other hand, adding too many network layers results in larger models that can easily overfit and learn a suboptimal profiling model, thus reducing the possibility of fitting the existing leakages. One straightforward way to avoid this problem is by allowing larger models to be trained with

- G. Perin is with the Delft University of Technology, The Netherlands.
- L. Wu is with the Delft University of Technology, The Netherlands.
- S. Picek is with Radboud University, The Netherlands.

regularization, restricting the model’s capacity during training. Dropout, weight decay, and data augmentation are well-known methods for regularization, but their indirect influence on the attack performance adds newly introduced hyperparameters to the tuning process. Alternatively, early stopping is a very efficient regularization mechanism that monitors a validation metric and saves model parameters (weights and biases) when the training reaches the best generalization moment.

An efficient early stopping implementation in profiling SCA requires monitoring the most appropriate metric. Reducing cross-entropy loss has been widely considered as the main training objective [11]. This is especially advantageous when security evaluations follow worst-case assumptions where the learning of an optimal model, which has good generalization, should also provide the smallest possible validation loss value. Unfortunately, in profiling SCA, collected leakages are normally extremely noisy because of environmental noise and implemented countermeasures. This usually leads the profiling process to end up in a suboptimal model. In those cases, as also empirically demonstrated in [12], validation loss (and accuracy) are inconsistent with SCA performance (i.e., key recovery) when the model is trained on protected datasets and, sometimes, overfits. Although the model can be optimized through gradient descent by minimizing generic loss functions (such as categorical cross-entropy or negative log-likelihood), the calculation of guessing entropy (GE) from a set of validation traces is consistent and highly informative concerning the profiling model generalization in SCA. The main reason for that is because GE measures the summation likelihood for each possible key guess over a set of traces instead of assessing the individual probabilities of expected classes only, as is the case of machine learning metrics. Therefore, the application of empirical GE as an early stopping metric tends to be reliable to assess model generalization during training. However, empirical GE provides significant overheads depending on the validation set size. If early stopping is adopted with hyperparameter tuning, the process becomes very slow and, in some cases, impractical.

Contributions: To address the unsolved problem of having a highly efficient early stopping metric for profiling SCA, we propose a fast guessing entropy calculation by simply reducing the number of validation traces when accessing model generalization during training. By doing so, we show that the trained models do not suffer in performance, but the training process becomes significantly faster, allowing more detailed tuning. Our fast GE method (denoted FGE) is especially important when security evaluators relax adversaries’ assumptions and do not assume the knowledge of secret random masks. For this reason, deep learning-based profiling attacks tend to become more difficult, which then requires a larger number of model search attempts. We compare the FGE method with state-of-the-art metrics for early stopping and guessing entropy in a deep learning-based SCA context. We show that FGE estimation is highly competitive and provides superior results with a neglectable time overhead in all scenarios. With FGE, training with early stopping becomes faster, allowing hyperparameter tuning to deploy more search attempts and increasing the chances to select the model with higher performance.

2 BACKGROUND

In this section, we start by providing details about deep learning-based SCA and commonly used metrics. Afterward, we discuss the datasets we use in our experiments.

2.1 Deep learning-based SCA

Profiling SCAs consider the strongest adversary with access to a clone device running the target cryptographic algorithm. The adversary can query the clone device with any set of plaintext $P = (p_0, p_1, \dots, p_{N-1})$ and chosen keys $K = (K_0, K_1, \dots, K_{N-1})$, and measure side-channel traces $X = (x_0, x_1, \dots, x_{N-1})$. These traces (X_{prof}) are used for training the classification algorithm (i.e., to build a machine learning model). This phase is known as the training or profiling phase. During this phase, a validation set, X_{val} , containing V traces, is selected from the profiling set to validate the model. Next, the adversary obtains measurements from the target device, where traces (X_{attack}) are also captured with known input but unknown (secret) key. The previously trained model is then exploited to recover the secret key k^* used in the target device. This phase is known as the attack or test phase.

The template attack is the first introduced profiling approach in SCA [1]. This attack is also the best possible attack if sufficient (infinite) training traces are available [13]. Over the years, machine learning and deep learning algorithms have been shown to be more powerful in realistic scenarios, where noise and countermeasures further reduce the measurement quality [8], [9]. While the profiling attack assumes a more powerful attacker than a non-profiling one, it requires significantly fewer traces than direct attacks to break the target: sometimes, only one trace could be sufficient.

Profiling SCA considers different methods to build or learn the statistical parameters representing a profiling model $f(\theta)$. The template attack [1] assumes that side-channel leakages follow a multivariate Gaussian distribution. The profiling phase consists of computing statistical parameters for a Gaussian mixture model (θ is given by mean and covariance parameters). Thus, the model is built for each possible hypothetical leakage class (e.g., all possible Hamming weight values of a byte). In the attack phase, the adversary computes the probability that a new side-channel measurement (under attack) belongs to a specific class by using the computed probability density function from the approximate statistics.

In the case of machine learning (including deep learning), the statistical parameters θ (e.g., weights and biases in the case of neural networks) are learned from profiling traces X_{prof} during the training phase. The deep neural network can skip feature selection from X_{prof} , which is an advantage over classic machine learning techniques and template attacks [14].

In the attack phase, the adversary obtains a probability v_k that the set of attack traces X_{attack} process the key byte $k \in K$, according to:

$$v_k = \sum_{i=0}^{Q-1} \log p[l(d_i, k)|x_i], \quad (1)$$

where $l(d_i, k)$ is the leakage function computed from public information d_i and the key hypothesis k . In our case, as we

attack AES implementations running encryption executions, the leakage function is given by $l(d_i, k) = S\text{-box}(d_i \oplus k)$ for the Identity leakage model or $l(d_i, k) = HW(S\text{-box}(d_i \oplus k))$ for the Hamming weight leakage model. The public value d_i is the corresponding plaintext byte. The recovered (guessed) key byte k from X_{attack} is then obtained as:

$$k = \underset{k \in K}{\operatorname{argmax}}([v_k]). \quad (2)$$

If the model is good (i.e., it learned the leakage), then the recovered key is k^* or at least k^* is among the best guesses.

2.2 Metrics

The training process has the minimization of the selected loss function as the main goal. In this paper, we consider the *categorical cross-entropy* (CCE) as the loss function. As demonstrated in [12], due to the imbalanced dataset problem, validation loss function values (including CCE) can be inconsistent with SCA metrics, which is also the case of SCA-based loss functions, as already proposed in [15], [16]. Therefore, we must select a more efficient validation metric to assess the model's performance for SCA.

Metrics like guessing entropy (GE) [17] are commonly used by an adversary to estimate the required effort to obtain the key. A side-channel attack outputs a key guessing vector $\mathbf{g} = [g_1, g_2, \dots, g_{|\mathcal{K}|}]$ in decreasing order of probability, i.e., g_1 represents the most likely key candidate and $g_{|\mathcal{K}|}$ the least likely key candidate. Guessing entropy is the average position of k^* in \mathbf{g} . Commonly, the averaged value is calculated over multiple independent experiments to obtain statistically significant results. In this paper, this GE method is called empirical GE, and it is evaluated on a set of V validation traces, where the results of multiple key rank executions are averaged and performed on a partition Q from V .

2.3 Datasets

We consider three datasets commonly used in research on deep learning-based SCA.

2.3.1 ASCAD

We evaluate ASCAD datasets [18] that contain side-channel measurements collected from the first-order protected software implementations of AES-128 running on an 8-bit AVR microcontroller [19]. There are two versions of the ASCAD dataset. The first version, ASCADf, has a fixed key and 60 000 traces in total. We split the dataset into 50 000, 5 000, and 5 000 for profiling, validation, and attack sets, respectively. The second version of the ASCAD dataset, ASCADr, has fixed and random keys, and it consists of 300 000 traces. In this case, we consider 200 000 for profiling (with random keys), 10 000 for validation, and 10 000 for the attack set. Both validation and attack sets have a fixed key. For both versions, we attack the third key byte (which is the first masked byte) by using the trimmed intervals already extracted and released in [18]. Thus, we use a pre-selected window of 700 features for ASCADf, while for ASCADr, the window size equals 1 400 features. For all experiments, the datasets are labeled according to the leakage model from the third S-box output byte in the first AES encryption round, i.e., $S\text{-box}(p_i \oplus k_i)$ and $HW(S\text{-box}(p_i \oplus k_i))$ for the Identity and Hamming weight leakage models, respectively.

2.3.2 CHES CTF 2018

This AES dataset was released as part of the Capture-the-Flag (CTF) ¹ competition in the Cryptographic Hardware and Embedded Systems (CHES) workshop in 2018. Four sets of 10 000 traces, featuring encryption operations of a first-order masked software implementation, were released for profiling purposes. These four sets were measured from four different STM32 platforms, namely A, B, C, and D. Additional two sets of 1 000 traces were released as attack traces from devices C and D. In our experiments, we consider the three first sets A, B, and C, containing random keys and random inputs, as a set of 30 000 profiling traces. The set from device D, containing the fixed key, is then used as attack and validation sets. As side-channel measurements from CHES CTF contain 650 000 samples points per trace, we performed a window resampling on the traces and concatenated two intervals representing the mask processing before encryption and target intermediate operation, i.e., the S-Box in the first encryption round. The resulting dataset contains 4 000 sample points. Both trace intervals are selected through a visual trace inspection. Note that for this dataset, source code and secret mask shares are not provided.

3 RELATED WORKS

Optimizing performance in DL-based profiling SCA has received significant attention in recent years. Due to the expensive trial-and-error cost in the profiling phase, enhancing performance in DL-based profiling SCA is a challenging task. In recent years, the SCA community considered two main alternatives to improve the attack efficiency: (1) by defining small neural network models that are faster to train and easier to tune [8], [9] and (2) by reducing the number of the required profiling traces during training [20]. Both solutions can have severe impacts on attack or generalization performance. The first approach may result in models that underfit for more noisy leakages or leakages obtained from other devices (portability problem [21]). The second alternative speeds up the process; still, it may result in limited learnability due to the eventually low number of profiling traces. Naturally, to reach small neural networks models, one needs to use appropriate methodology. Zaid et al. [8] and Wouters et al. [22] worked on designing methodologies for finding efficient neural network architectures. Rijdsdijk et al. [9] and Wu et al. [23] investigated advanced hyperparameter tuning techniques like reinforcement learning and Bayesian optimization, respectively. Perin et al. showed that even a random search could find very successful neural network models [24].

Besides the methods mentioned above, a third alternative uses efficient and reliable validation metrics to evaluate training and, consequently, implement faster hyperparameter tuning (which can provide faster convergence) with larger models and larger profiling sets. Empirical GE (described in Section 2.1) can be very expensive to compute with larger validation sets, especially if used during training to detect the best training epoch. In a recent publication, Zhang et al. [15] proposed Guessing Entropy Estimation

1. <https://chesctf.riscure.com/2018/content?show=training>

Algorithm (GEEA) to reduce the computational limitation cost of empirical GE for the full attacked key scenarios, which computes faster than empirical GE calculation on separate key bytes. Indeed, empirical GE executes multiple key rank executions over multiple partitions of the dataset V , each partition containing Q measurements. GEEA, on the other hand, only requires one execution over the Q measurements.

Let us consider $s(k_g, x_i, d_i)$ as a score indicating the probability that a measurement x_i process key k_g for a input (i.e., plaintext) d_i . The GEEA first requires the calculation of pairwise subtractions of scores concerning the correct key, resulting in mean and variance for each key guess $k_g \in K$ as follows:

$$\mu_{k_g} = \frac{1}{Q} \sum_{i=0}^{Q-1} [s(k_g, x_i, d_i) - s(k_c, x_i, d_i)] \quad (3)$$

$$\sigma_{k_g} = \sqrt{\frac{1}{Q} \sum_{i=0}^{Q-1} [s(k_g, x_i, d_i) - s(k_c, x_i, d_i) - \mu_{k_g}]^2}, \quad (4)$$

where k_c is the correct key. Then, guessing entropy value is obtained as:

$$GE = 1 + \sum_{k=0, k_g \neq k_c}^{|K|} \Phi\left(\frac{\sqrt{Q}\mu_k}{\sigma_k}\right), \quad (5)$$

where $\Phi()$ is the cumulative density function of a normal Gaussian distribution $\mathcal{N}(0, 1)$.

Alternative solutions were proposed as new validation metrics for early stopping, stopping training sooner, and speeding up the process. In [25], the authors considered mutual information approach between the output probabilities and validation labels to monitor the best epoch during training. The work of [26] monitors the epoch when the training achieves the minimal difference between the number of profiling and validation traces that are required to achieve 90% of success rate. The authors proposed a routine to abort training if this difference increases after reaching its minimum value. In our work, we also consider the mutual information metric for comparison. The method proposed in [26] is not considered in our comparative analysis as it is directly adapted to datasets with fixed keys in the profiling set, which is not the case of ASCADr. The method requires estimating the number of traces to reach a success rate of 90%, which implies obtaining the evolution of success rate concerning the number of validation traces. This means that the success rate is computed Q times for each epoch, adding significant time overhead to the process.

As we can see, none of the mentioned approaches compute GE directly from the validation traces at the end of each training epoch. GEEA was proposed as a fast and more stable GE estimation, but it is not suggested to be used during network training. On the other hand, although GE can be a potential metric candidate, its computation could be very slow if more validation traces are considered (which is required for GE stability), finally providing significant overheads to the training process. Therefore, the SCA community did not consider directly applying GE (including GEEA) as the early stopping metric, especially

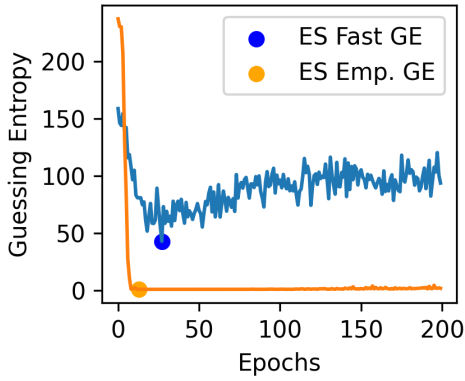
in the hyperparameter search processes. This work shows that significantly reducing the number of validation traces for GE estimation during training is a reliable and efficient metric for early stopping, benefiting hyperparameter tuning optimization.

4 FAST GE FOR EARLY STOPPING

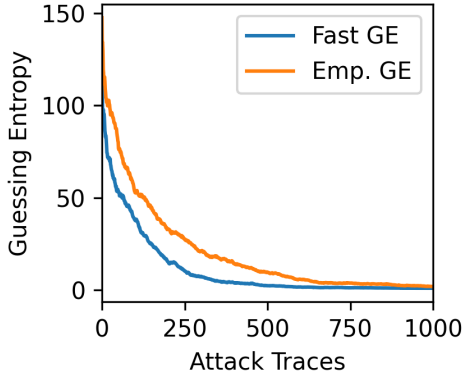
Running hyperparameter search without pre-selecting efficient ranges for each hyperparameter may fail to find powerful attack models. A solution could be searching for small models with restricted search ranges, as proposed in [8], [22] or by setting the objective of the search as being a small model, as proposed by Rijdsdijk et al. [9]. Small models are usually self-regularized, but they still suffer from limited fitting capacity, which is particularly problematic for noisy and protected targets. An alternative is to allow larger models and add regularization to prevent overfitting [27]. Although regularization tends to improve model generalization, regularized models with increased size require more training epochs, reducing the efficiency in a hyperparameter search process. As the number of training epochs is a critical hyperparameter to be determined, early stopping may become a standard approach.

To allow efficient early stopping with GE, we propose a fast GE (FGE) calculation to reduce the empirical GE overheads. When used as an early stopping metric, FGE provides very small overheads to the training process, usually between 1.5% and 3.3%, while, e.g., empirical GE shows overheads between 18.59% and 28.19%, as reported in Section 5. Our idea consists in reducing the number of validation traces when computing GE for each processed epoch, which has multiple benefits in DL-based SCA. The pseudo-code showing how FGE is obtained is provided in Algorithm 1. As the algorithm shows, the main application of FGE is for the hyperparameter search process.

If the model converges, the attack is successful, and the GE for a small number of validation traces can also indicate the best epoch to stop training efficiently. Using large validation sets for the metric calculation may obscure the real performance of the model: a model that overfits may also slowly decrease guessing entropy to 1 after processing enough validation traces. In contrast, FGE is more sensitive to the model's performance change, thanks to its low usage of the validation traces. This situation is illustrated in Figure 1. As we can see in Figure 1a, when using $Q = 3000$ validation traces for empirical GE, the best epoch will be returned at the moment when GE is equal to 1. If the next epochs indicate a model that requires even fewer traces to succeed (which means better generalization), empirical GE will not capture that. On the other hand, using fewer traces allows us to obtain this convergence and, as a consequence, be able to recover the key with fewer attack traces, as shown in Figure 1b. Of course, the question here is: why would this be a problem if reaching GE of 1 allows an adversary to recover the key? We can observe two main problems in this scenario. First, empirical GE with more traces provides more overhead and limits the number of hyperparameter search attempts, preventing us from finding a model that eventually breaks the target (which is the example of model found in Figure 1). Second, from the current model, we



(a) GE vs epochs.



(b) GE vs attack traces.

Fig. 1: Fast GE vs Empirical GE (ES = Early Stopping).

would select trained parameters before it reaches its best attack performance or generalization capacity, which can also indicate overfitting on the validation set, possibly opening issues in portability scenarios (when the device used for profiling is different from the device used for attack [21]). If a model generalizes, then GE will eventually decrease, and FGE should show this behavior too.

5 EXPERIMENTAL RESULTS

This section provides experimental results for 1) machine learning models obtained through a random search, 2) hyperparameter tuning for different validation metrics, and 3) state-of-the-art model and different validation metrics.

5.1 Hyperparameter Search Ranges

In this work, we only consider convolutional neural networks (CNNs) as they contain many hyperparameters to tune and, therefore, it becomes more challenging to find good hyperparameter combinations than, e.g., multilayer perceptrons.

Convolutional neural networks commonly consist of three types of layers: convolutional layers, pooling layers, and fully connected layers. The convolution layer computes the output of neurons connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. Pooling decrease the number of extracted features

by performing a down-sampling operation along the spatial dimensions. Finally, the fully connected layer computes the hidden activations or the class scores.

Table 1 provides the selected ranges for the hyperparameter tuning processes. These selected ranges result in a search space containing 2.7×10^9 possible combinations. As we can see, we allow CNNs to contain up to eight hidden layers, combining convolution and dense layers. A pooling layer always follows each convolution layer. As the ASCADf and ASCADr datasets contain 50 000 and 200 000 profiling traces, respectively, larger models would tend to overfit.

TABLE 1: Hyperparameter search space for CNNs (c in convolution filters indicates the convolution layer index).

Hyperparameter	Ranges		
	Min	Max	Step
Batch Size	100	1 000	100
Convolution Layers	1	4	1
Convolution Filters	$2 \times 2^{c-1}$	$16 \times 2^{c-1}$	2
Kernel Size	4	20	1
Stride	1	4	1
Pooling Size	1	4	1
Pooling Stride	1	4	1
Dense layers	1	4	1
	Options		
Neurons	10, 20, 30, 40, 50, 100, 200, 300, 400, 500		
Activation function	ReLU, ELU, or SELU		
Learning Rate	5e-3, 1e-3, 5e-4, 1e-4, 5e-5, 1e-5		
Pooling Type	Average, Max		
Weight Initializer	He, Random Uniform, Glorot Uniform		
Optimizer	Adam, RMSprop		

5.2 Random Hyperparameter Search with Different Validation Metrics

We compare different early stopping metrics in a random hyperparameter search process for the two ASCAD datasets and different leakage models. The results of the CHES CTF dataset are only provided with Hamming weight leakage model. Each randomly selected CNN is trained for 200 epochs, and we save the trained weights at the end of each epoch. At the end of the training, each early stopping metric indicates the best training epoch, and we restore the trained weights from that epoch. Then, as the training is finished, we compute GE for the attack set containing a larger number of traces. Note that 200 epochs is a relatively small number for training epochs, and, as shown in this section, stopping the training after 200 epochs may also deliver good results for some cases.

Table 2 gives the number of validation traces V considered for each early stopping metric, while the partition amount Q is the number of the traces used to calculate each specific metric. For instance, GE is the average of multiple key rank executions over Q traces, which are randomly selected from a larger set V for each key rank executions. This way, we set V greater than Q so that sampling each data in Q preserves certain randomness. By doing so, the obtained results would indicate a better generalization capacity of models. For mutual information, we apply V validation traces. FGE estimation considers only 50 traces for Q and 500 traces for V . We tested other values for Q and V , from

Algorithm 1 Hyperparameter Search with Early Stopping and Fast Guessing Entropy (FGE).

```

1: Set  $\Theta$  as the set of models
2: for new search attempt  $S$  do
3:   Generate new hyperparameter set  $\mathcal{H}$ 
4:   Initialize model parameters  $\theta$ 
5:   Select a small validation set  $X_{val-fast}$ 
6:   for Epoch  $E$  in Epochs do
7:     Train neural network model  $\mathcal{F}_{\mathcal{H}}(\theta, X_{prof}, Y_{prof})$ 
8:     Compute FGE:  $GE_{fast}[E] = GE(\theta, X_{val-fast}, Y_{val-fast})$ 
9:     Save model parameters  $\theta[E]$  at epoch  $E$ 
10:  end for
11:  Retrieve model parameters from best epoch:  $\Theta[S] = \theta[\text{argmin}(GE_{fast})]$ 
12:  Select a full validation set  $X_{val}$ 
13:  Compute  $GE(\theta_{best}, X_{val}, Y_{val})$  and the corresponding number of validation traces to reach  $GE = 1$ ,  $N_{GE1}[S]$ .
14: end for
15: Return best model:  $\theta_{best\_model} = \Theta[\text{argmin}(N_{GE1})]$  (model that requires minimum number of validation traces to reach  $GE = 1$ ).

```

20 up to 200, and 50 was the minimum value for Q and V that still preserves the best results for FGE.

TABLE 2: Number of validation traces for each early stopping method.

	ASCADf		ASCADr		CHES CTF	
	V	Q	V	Q	V	Q
Fast GE	500	50	500	50	500	50
Emp. GE	5000	3000	10000	5000	5000	3000
GEEA	5000	3000	10000	5000	5000	3000
MI	5000	5000	10000	10000	5000	5000

We execute 500 searches for each dataset, considering the Hamming weight and Identity leakage models. Table 3 provides the average time overhead in percentage for each considered metric. As we can see, the FGE estimation provides a maximum of 3.35% overhead among the four considered scenarios. For the ASCADr dataset, the overhead is only 1.19% and 1.49%, which can be considered negligible for the training time compared with its counterparts. As expected, empirical GE and GEEA methods provide the largest overheads, although GEEA is faster than empirical GE. The mutual information method provides the second-best results, which is related to the more straightforward calculation than guessing entropy.

TABLE 3: Average time overhead of different early stopping methods.

	ASCADf		ASCADr		CHES CTF
	HW	Identity	HW	Identity	HW
Fast GE	2.66%	3.35%	1.19%	1.49%	2.74%
Emp. GE	20.70%	28.19%	18.59%	24.21%	20.61%
GEEA	11.48%	23.95%	9.31%	20.17%	13.56%
MI	9.28%	7.46%	7.81%	6.30%	11.28%

Table 4 provides the % that each metric can select a generalizing model with early stopping (model that reaches $GE=1$ in the attack phase, which is indicated by line 13 in Algorithm 1) from the random search. Together with

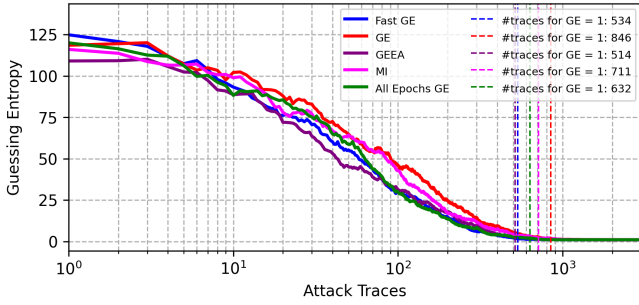
GEEA, the fast GE is a highly efficient metric (top two performance in all considered scenarios). Most importantly, we successfully verify that FGE is always superior to the situation where no early stopping is used (200 epochs in the table) and with neglectable overhead. For the case of the Identity leakage models, FGE shows the best results.

TABLE 4: % of times a generalizing DNN was selected from each metric and from the training with all 200 epochs.

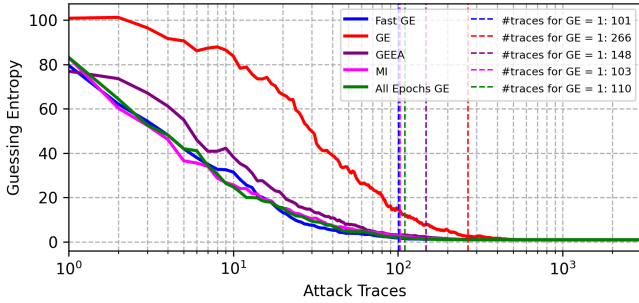
	ASCADf		ASCADr		CHES CTF
	HW	Identity	HW	Identity	HW
Fast GE	56.52%	43.46%	49.63%	34.13%	21.85%
Emp. GE	59.66%	43.16%	50.00%	33.23%	20.79%
GEEA	59.66%	43.46%	54.34%	29.30%	21.18%
MI	50.74%	37.38%	43.84%	33.53%	19.78%
200 epochs	49.75%	40.42%	45.83%	32.62%	15.06%

Figure 2 shows results for the ASCADf dataset. When side-channel traces are labeled according to the Hamming weight leakage model, the correct key is recovered with 514 traces for GEEA metric and 534 traces (the second best) with FGE early stopping metric. In the case of the Identity leakage model, the best results are achieved for the FGE metric, where 101 attack traces are needed to achieve GE equal to 1, which is aligned with state-of-the-art results [8], [9], [22]. The good performing results from the mutual information metric and the GE obtained with 200 epochs indicate the effectiveness of early-stopping metrics in preventing the best model from overfitting. Again, we confirm that FGE is highly competitive in both leakage models and requires $10\times$ fewer validation traces.

For ASCADr dataset, results for FGE are also very promising, as shown in Figure 3. For the Hamming weight leakage model, FGE provides the best results, followed by mutual information metric. In the case of the Identity leakage model, the best result is obtained with all 200 epochs, showing that this number of epochs is appropriate for this best model found through random search. The best results are obtained with the FGE metric when early stopping is considered.

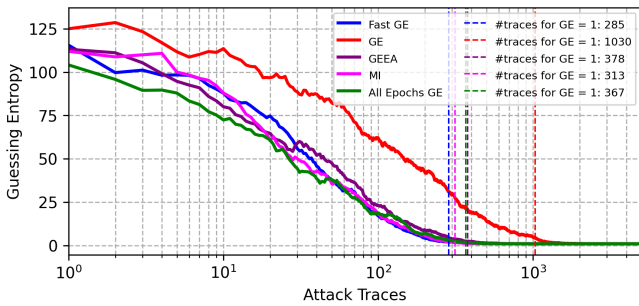


(a) Hamming weight leakage model.

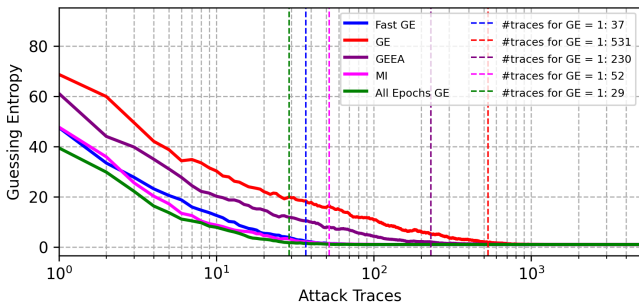


(b) Identity leakage model.

Fig. 2: GE results from best models selected from different early stopping metrics for the ASCADf dataset.



(a) Hamming weight leakage model.



(b) Identity leakage model.

Fig. 3: GE results from best models selected from different early stopping metrics for the ASCADr dataset.

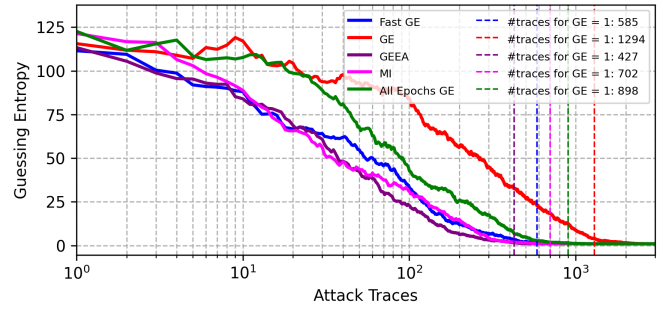


Fig. 4: GE results from best models selected from different early stopping metrics for the CHES CTF 2018 dataset (Hamming weight leakage model).

Figure 4 provides results for the CHES CTF dataset. FGE metric provides second-best results after GEEA. Results for the CHES CTF dataset are only shown for the Hamming weight leakage model, as this dataset provides bad results with the Identity leakage model, as discussed in [9].

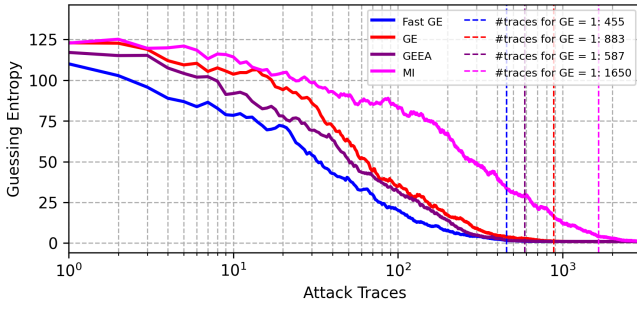
Furthermore, the performance of the best models selected from empirical GE as an early stopping metric provided less efficient results. As already mentioned in [15], empirical GE requires a very large validation set and a more stable GE estimation can be obtained with the selection of larger validation sets. Of course, using larger validation sets provides an estimation of model generalization. This is especially important for models that provide suboptimal performance and require more traces to show guessing entropy reduction for the correct key. However, computing GE for this large number of traces is undesirable as an early stopping metric due to significant time overhead.

5.3 Hyperparameter Tuning with Different Validation Metrics

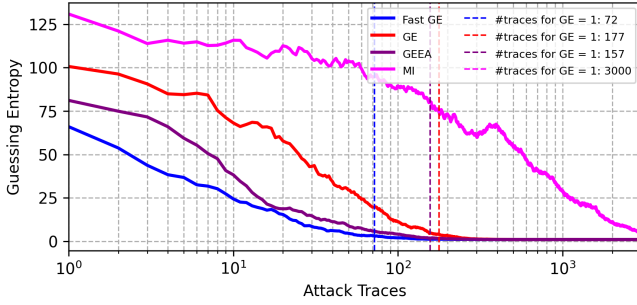
This section analyzes how the evaluated early stopping metrics perform with Bayesian optimization (BO) for hyperparameter search [23]. For that, we consider the open-source `BayesianOptimization` method provided in `keras-tuner` [28] Python package. We run BO for 100 searches with ASCAD datasets and the Hamming weight and Identity leakage models. We repeat each search process five times for each different early stopping metric. The guessing entropy results without early-stopping (“all epochs” labels in figures from the previous section) are omitted because `keras-tuner` inherently implements early-stopping and, for this reason, it is not possible to select the best model by ignoring early-stopping. The results reported in this section are extracted from the best-found model out of the five search attempts.

Results from BO for the ASCADf dataset are shown in Figure 5. The best results are obtained from FGE for both Hamming weight and Identity leakage models. In particular, for the Identity leakage model, as shown in Figure 5b, the best found model achieves GE equal to 1 with less than half of the attack traces needed for GEEA. In these experiments, mutual information provides less efficient results.

Figure 6 provides BO results for the ASCADr dataset. For the Hamming leakage model, GEEA and FGE provide the best results. For the Identity leakage model, results for FGE



(a) Hamming weight leakage model.



(b) Identity leakage model.

Fig. 5: GE results from best models found with BO with different early stopping metrics for the ASCADf dataset.

are superior, and only 60 attack traces are required for key byte recovery, while empirical GE requires $10\times$ more attack traces to succeed. Again, the mutual information metric delivers the worst results.

Running hyperparameter tuning with Bayesian optimization for the CHES CTF dataset and the Hamming weight leakage model, the results obtained with FGE are significantly better compared to other validation metrics, as shown in Figure 7. As we can see, FGE returns the best model that reaches GE equal to 1 in the attack phase with only 232 traces, while other metrics always significant more attack traces.

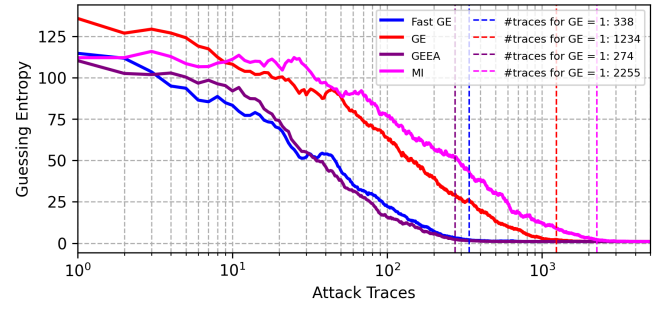
5.4 State-Of-The-Art Models with Different Validation Metrics

The works of [8], [9], [22] proposed hyperparameter tuning for ASCADf dataset and their models reported state-of-the-art results. In this section, we also verify how FGE can improve the performance of those best models even more. This way, we provide attack results when applying early stopping to three different CNN architectures. As the results for these CNN were reported for the Identity leakage model, we only consider that scenario in our analysis.

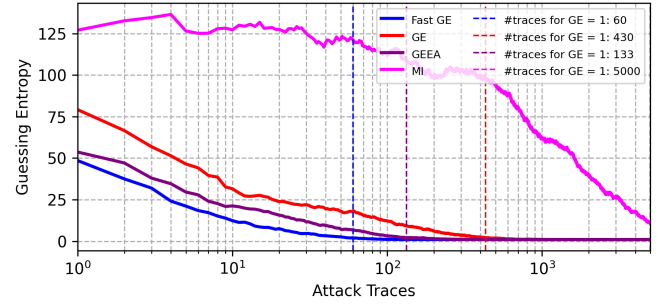
As shown in Figure 8, for CNN models from [8] and [9], our FGE metric provides the best results. Results for CNN model from [22] also put FGE among the best-performing metrics.

6 CONCLUSIONS AND FUTURE WORK

Profiling attacks are important during security evaluations because evaluators can determine if the device leaks information with high assurance. This is especially possible



(a) Hamming weight leakage model.



(b) Identity leakage model.

Fig. 6: GE results from best models found with BO with different early stopping metrics for the ASCADr dataset.

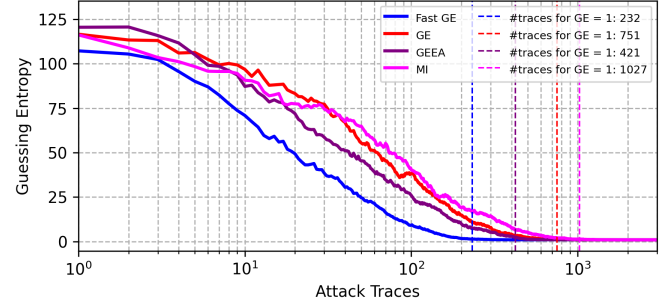


Fig. 7: GE results from best models found with BO with different early stopping metrics for the CHES CTF dataset (Hamming weight leakage model).

because assumptions during a profiling analysis consider that the target faces an adversary that can learn existing side-channel leakages in a supervised learning setting.

In a recent publication [29], Bronchain et al. showed through the lens of perceived information (PI) [30] how different profiling methods perform against protected cryptographic implementations. Their analysis allows security evaluators to conclude about the target's leakages with the worst-case security. For that, the evaluator assumes that the adversary has knowledge of all intermediate secret shares during profiling as well as the source code. Consequently, such an evaluation provides conditions to implement optimal profiling models, where assumptions about the target (e.g., leakage model) contain as few errors as possible. Also, the evaluator can build a profiling model with a sufficient number of traces, thus minimizing estimation errors.

The case of deep learning for profiling SCA brings a new

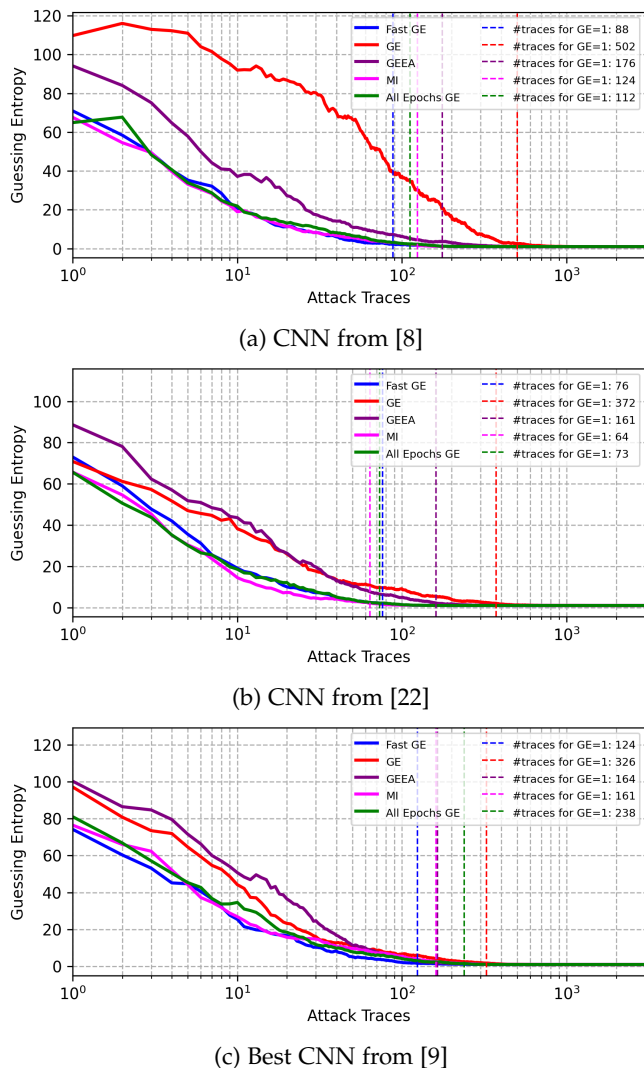


Fig. 8: Performance of different validation metrics on state-of-the-art CNN architectures, ASCADf dataset.

perspective over profiling attacks. The main reason for that is related to the ability of a deep neural network to perform efficiently without feature selection. In practice, this means that the attacked interval contains several low SNR points of interest, and the selection of most leaky points of interest with a high SNR becomes more challenging. The advantages for security evaluations come from the fact that neural networks as profiling models can learn existing leakages even without feature selection and, in practice, deliver close to optimal results. The results for CNN architectures from Figure 8 is an example of this case. Of course, to reach an optimal deep learning profiling model, costly hyperparameter tuning needs to be implemented, especially for more protected targets.

Therefore, to reach optimal deep learning models without worst-case security assumptions, hyperparameter tuning needs to be as efficient as possible. For that, assessing the model generalization during training becomes crucial, requiring fast and efficient validation metrics. We propose using a fast GE metric that requires significantly fewer validation traces in the GE calculation. Our results indicate

that FGE as a validation metric delivers efficient and competitive early stopping results. Our technique is validated in different scenarios and shows good results with neglectable time overheads. Thus, we consider FGE as the method of choice for practical deep learning-based SCA hyperparameter tuning.

As future works, we will explore the efficiency of different validation metrics in portability settings and with different countermeasures in future work. Additionally, as this work contains results for convolutional neural networks only, it would be interesting to assess FGE performance with architectures like multilayer perceptrons and residual neural networks.

REFERENCES

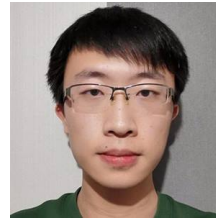
- [1] S. Chari, J. R. Rao, and P. Rohatgi, "Template attacks," in *Cryptographic Hardware and Embedded Systems - CHES 2002*. Springer Berlin Heidelberg, 2003, pp. 13–28.
- [2] W. Schindler, K. Lemke, and C. Paar, "A stochastic model for differential side channel cryptanalysis," in *Cryptographic Hardware and Embedded Systems - CHES 2005, 7th International Workshop, Edinburgh, UK, August 29 - September 1, 2005, Proceedings*, ser. Lecture Notes in Computer Science, J. R. Rao and B. Sunar, Eds., vol. 3659. Springer, 2005, pp. 30–46. [Online]. Available: https://doi.org/10.1007/11545262_3
- [3] P. C. Kocher, J. Jaffe, and B. Jun, "Differential power analysis," in *Advances in Cryptology - CRYPTO '99, 19th Annual International Cryptology Conference, Santa Barbara, California, USA, August 15-19, 1999, Proceedings*, ser. Lecture Notes in Computer Science, M. J. Wiener, Ed., vol. 1666. Springer, 1999, pp. 388–397. [Online]. Available: https://doi.org/10.1007/3-540-48405-1_25
- [4] E. Brier, C. Clavier, and F. Olivier, "Correlation power analysis with a leakage model," in *Cryptographic Hardware and Embedded Systems - CHES 2004: 6th International Workshop Cambridge, MA, USA, August 11-13, 2004. Proceedings*, ser. Lecture Notes in Computer Science, M. Joye and J. Quisquater, Eds., vol. 3156. Springer, 2004, pp. 16–29. [Online]. Available: https://doi.org/10.1007/978-3-540-28632-5_2
- [5] B. Gierlichs, L. Batina, P. Tuyls, and B. Preneel, "Mutual information analysis," in *Cryptographic Hardware and Embedded Systems - CHES 2008, 10th International Workshop, Washington, D.C., USA, August 10-13, 2008. Proceedings*, ser. Lecture Notes in Computer Science, E. Oswald and P. Rohatgi, Eds., vol. 5154. Springer, 2008, pp. 426–442. [Online]. Available: https://doi.org/10.1007/978-3-540-85053-3_27
- [6] V. Banciu, E. Oswald, and C. Whittall, "Reliable information extraction for single trace attacks," in *Proceedings of the 2015 Design, Automation & Test in Europe Conference & Exhibition, DATE 2015, Grenoble, France, March 9-13, 2015*, W. Nebel and D. Atienza, Eds. ACM, 2015, pp. 133–138. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2755783>
- [7] L. Lerman, G. Bontempi, and O. Markowitch, "A machine learning approach against a masked AES - reaching the limit of side-channel attacks with a learning model," *J. Cryptogr. Eng.*, vol. 5, no. 2, pp. 123–139, 2015. [Online]. Available: <https://doi.org/10.1007/s13389-014-0089-3>
- [8] G. Zaid, L. Bossuet, A. Habrard, and A. Venelli, "Methodology for efficient CNN architectures in profiling attacks," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 1, pp. 1–36, 2020.
- [9] J. Rijdsdijk, L. Wu, G. Perin, and S. Picek, "Reinforcement learning for hyperparameter tuning in deep learning-based side-channel analysis," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2021, no. 3, pp. 677–707, 2021.
- [10] E. Cagli, C. Dumas, and E. Prouff, "Convolutional neural networks with data augmentation against jitter-based countermeasures - profiling attacks without pre-processing," in *Cryptographic Hardware and Embedded Systems - CHES 2017 - 19th International Conference, Taipei, Taiwan, September 25-28, 2017, Proceedings*, ser. Lecture Notes in Computer Science, W. Fischer and N. Homma, Eds., vol. 10529. Springer, 2017, pp. 45–68. [Online]. Available: https://doi.org/10.1007/978-3-319-66787-4_3

- [11] L. Masure, C. Dumas, and E. Prouff, "A comprehensive study of deep learning for side-channel analysis," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 1, pp. 348–375, 2020. [Online]. Available: <https://doi.org/10.13154/tches.v2020.i1.348-375>
- [12] S. Picek, A. Heuser, A. Jovic, S. Bhasin, and F. Regazzoni, "The curse of class imbalance and conflicting metrics with machine learning for side-channel evaluations," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2019, no. 1, pp. 209–237, 2019.
- [13] L. Lerman, R. Poussier, G. Bontempi, O. Markowitch, and F.-X. Standaert, "Template Attacks vs. Machine Learning Revisited (and the Curse of Dimensionality in Side-Channel Analysis)," in *Lecture Notes in Computer Science*, 2015, vol. 9064, no. ML, pp. 20–33.
- [14] X. Lu, C. Zhang, P. Cao, D. Gu, and H. Lu, "Pay attention to raw traces: A deep learning architecture for end-to-end profiling attacks," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2021, no. 3, p. 235–274, Jul. 2021. [Online]. Available: <https://tches.iacr.org/index.php/TCHES/article/view/8974>
- [15] J. Zhang, M. Zheng, J. Nan, H. Hu, and N. Yu, "A novel evaluation metric for deep learning-based side channel analysis and its extended application to imbalanced data," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 3, pp. 73–96, 2020.
- [16] G. Zaid, L. Bossuet, F. Dassance, A. Habrard, and A. Venelli, "Ranking loss: Maximizing the success rate in deep learning side-channel analysis," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2021, no. 1, pp. 25–55, 2021.
- [17] F. X. Standaert, T. G. Malkin, and M. Yung, "A unified framework for the analysis of side-channel key recovery attacks," *Lecture Notes in Computer Science*, vol. 5479 LNCS, pp. 443–461, 2009.
- [18] "ASCAD GitHub Repository," Website, 2018, <https://github.com/ANSSI-FR/ASCAD>.
- [19] R. Benadjila, E. Prouff, R. Strullu, E. Cagli, and C. Dumas, "Deep learning for side-channel analysis and introduction to ASCAD database," *J. Cryptographic Engineering*, vol. 10, no. 2, pp. 163–188, 2020.
- [20] S. Picek, A. Heuser, G. Perin, and S. Guilley, "Profiling side-channel analysis in the efficient attacker framework," *Cryptology ePrint Archive*, Report 2019/168, 2019.
- [21] S. Bhasin, A. Chattopadhyay, A. Heuser, D. Jap, S. Picek, and R. R. Shrivastwa, "Mind the portability: A warriors guide through realistic profiled side-channel analysis," in *27th NDSS*, 2020.
- [22] L. Wouters, V. Arribas, B. Gierlichs, and B. Preneel, "Revisiting a methodology for efficient CNN architectures in profiling attacks," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 3, pp. 147–168, 2020.
- [23] L. Wu, G. Perin, and S. Picek, "I choose you: Automated hyperparameter tuning for deep learning-based side-channel analysis," *Cryptology ePrint Archive*, Report 2020/1293, 2020.
- [24] G. Perin, L. Chmielewski, and S. Picek, "Strength in numbers: Improving generalization with ensembles in machine learning-based profiled side-channel analysis," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2020, no. 4, pp. 337–364, Aug. 2020. [Online]. Available: <https://tches.iacr.org/index.php/TCHES/article/view/8686>
- [25] G. Perin, I. Buhan, and S. Picek, "Learning when to stop: A mutual information approach to prevent overfitting in profiled side-channel analysis," ser. LNCS, vol. 12910. Springer, 2021, pp. 53–81.
- [26] D. Robissout, G. Zaid, B. Colombier, L. Bossuet, and A. Habrard, "Online performance evaluation of deep learning networks for profiled side-channel analysis," ser. Lecture Notes in Computer Science, vol. 12244. Springer, 2020, pp. 200–218.
- [27] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 56, pp. 1929–1958, 2014.
- [28] T. O'Malley, E. Bursztein, J. Long, F. Chollet, H. Jin, L. Invernizzi et al., "Kerastuner," <https://github.com/keras-team/keras-tuner>, 2019.
- [29] O. Bronchain, F. Durvaux, L. Masure, and F.-X. Standaert, "Efficient profiled side-channel analysis of masked implementations, extended," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 574–584, 2022.
- [30] O. Bronchain, J. M. Hendrickx, C. Massart, A. Olshevsky, and F. Standaert, "Leakage certification revisited: Bounding model errors in side-channel security evaluations," in *Advances in Cryptology - CRYPTO 2019 - 39th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 18-22, 2019, Proceedings*,

Part I, ser. Lecture Notes in Computer Science, A. Boldyreva and D. Micciancio, Eds., vol. 11692. Springer, 2019, pp. 713–737. [Online]. Available: https://doi.org/10.1007/978-3-030-26948-7_25



Guilherme Perin is a postdoctoral researcher at the Delft University of Technology. He graduated in Electrical Engineering (2008) and has Master in Informatics (2011) by the Federal University of Santa Maria. In 2014, he received his PhD in Microelectronics and Automated Systems at University of Montpellier. His research areas include hardware security, cryptography, optimization algorithms, and machine learning.



Lichao Wu is a PhD student in the cybersecurity research group at the Delft University of Technology. After obtaining a bachelor's degree at Northwestern Polytechnical University (2015), Wu received his master's degree in Microelectronic at the Delft University of Technology in 2017. His main research interests are at the intersection of implementation attacks, cryptography, and machine learning.



Stjepan Picek is an associate professor at Radboud University, The Netherlands. He received his PhD in 2015, and from 2015 to 2017, he was a postdoctoral researcher at KU Leuven, Belgium and MIT, USA. From 2017 to 2021, he was an assistant professor at the Delft University of Technology, The Netherlands. His research interests include cryptography, machine learning, and evolutionary algorithms.