# Remove Some Noise: On Pre-processing of Side-channel Measurements with Autoencoders

Lichao Wu and Stjepan Picek

Delft University of Technology, The Netherlands
l.wu-4@tudelft.nl, s.picek@tudelft.nl

**Abstract.** In the profiled side-channel analysis, deep learning-based techniques proved to be very successful even when attacking targets protected with countermeasures. Still, this does not mean that countermeasures do not make the attacks more difficult or that deep learning attacks will always succeed. As such, to improve the performance of attacks, an intuitive solution is to remove the effect of countermeasures.

In this paper, we investigate whether we can consider certain types of countermeasures as noise and then use deep learning to remove that noise. We conduct a detailed analysis of four different types of noise and countermeasures either separately or combined and show that in all scenarios, denoising autoencoder improves the attack performance significantly.

## 1 Introduction

Side-channel analysis (SCA) is a threat exploiting weaknesses in physical implementations of cryptographic algorithms rather than the algorithms themselves [1]. During the execution of an algorithm, leakages like electromagnetic (EM) radiation [2] or power dissipation [3] can happen. With an unprotected target, an attacker can record the leakages and then retrieve the secret information by analyzing those leakages. Side-channel analysis can be divided into 1) direct attacks like single power analysis (SPA) and differential power analysis (DPA) [3], and 2) profiled attacks like template attack (TA) [4] and supervised machine learning-based attacks [5–8]. In recent years, machine learning-based approaches and especially deep learning-based approaches proved to be a powerful option when conducting profiled SCA. While such attack methods strongly threaten the security of cryptographic devices, there are still some serious limitations. More precisely, attack methods commonly rely on the correlation characteristics of the signal, i.e., signal patterns that are related to the data being processed. Once the correlation degrades, attacks become less effective and sometimes even useless [9]. To optimize the data correlation, common approaches are to delimit the Points of Interest (POI) [10] or improve the alignment of the leakage traces [11].

In some cases, the low signal-to-noise ratio of the leakage increases the difficulties of identifying these patterns. Additionally, there are various countermeasures in both hardware and software that make the attacks even more difficult. Such countermeasures can be divided into two categories: masking and hiding. Masking splits the sensitive intermediate values into different shares to decrease the key dependency [12,13]. Hiding, on the other hand, aims at reducing the side-channel information by adding randomness to the leakage signals. There are several approaches to hiding. For example, the direct addition of noise [14] or the design of dual-rail logic styles [15] are frequently considered options. Exploiting time-randomization is another alternative, e.g., Random Delay Interrupts (RDIs) [16] implemented in software and clock jitters from the hardware. Still, the countermeasures (especially the hiding ones) are not without weaknesses. No matter what hiding approach is used, we can treat the effects of countermeasures as noise because of their randomness. In other words, the ground truth of the traces always exists. If we can find a way to denoise the traces and recover the ground truth of the leakage, then the reconstructed traces could become more vulnerable to the SCA.

While considering the countermeasures as noise and then removing that noise sounds like an intuitive approach, this is not an easy problem. Indeed, the noise (both from the environment and countermeasures) is a part of a signal and those two components cannot be separated if we do not know their precise characterizations. While the signal can be relatively easily characterized, the noise component has a stochastic nature and thus, its characterization is not readily available. Even in the scenarios where countermeasures could be described with a function, we must take into account the environmental changes that will, in turn, change the noise characterization. Additionally, for realistic settings, we must also consider the portability and the differences among various devices [17]. Combining all these factors makes this problem very complex and to the best of our knowledge, there are no good (universal) approaches against it.

Common approaches to remove noise are using low-pass filters [18], conducting trace alignments [11], and various feature engineering methods [10,19]. More recently, the side-channel community used deep learning techniques that can conduct implicit feature selection and/or fight countermeasures [7,8,20]. While such techniques are useful, they are usually aimed either against a single source of noise or in cases when they can handle more sources of noise, they do not offer interpretability of results, i.e., it is not clear at what point noise removal stops and attack itself starts (or even if there is such a point). Finally, we emphasize that being able to reduce the noise comprehensively brings several advantages 1) understanding the attack techniques better, 2) understanding the noise better and consequently, (hopefully) being able to design stronger countermeasures, 3) ability to mount stronger/simpler attacks as there is no noise to consider.

In this paper, we propose a new method to remove the hiding countermeasure with a denoising autoencoder. Although the denoising autoencoder has been proved to be successful to remove the noise from an image [21], to the best of our knowledge, this technique has not been applied to the side-channel domain to

reduce the noise/countermeasures effect. In this paper, we demonstrate the effectiveness of a convolutional denoising autoencoder in dealing with different types of noise and countermeasures separately, i.e., white noise, desynchronization, RDIs, and clock jitters. Then, we make the problem more realistic by combining various types of noise and countermeasures with the traces and trying to denoise it with the same machine learning model. The results show that the denoising autoencoder is surprisingly efficient in removing the noise and countermeasures in all investigated situations.

## 1.1 Related Work

Analysis of the leakage traces in the profiled SCA scenario can be seen as a classification problem where the goal of an attacker is to classify leakage traces based on the related data (i.e., the encryption key). The most powerful attack from the information-theoretic point of view is the template attack (TA) [4]. Still, this attack can reach its full potential only if the attacker has an unbounded number of traces and the noise follows the Gaussian distribution [22]. More recently, various machine learning techniques emerged as preferred options for cases where 1) the number of traces is either limited or very large, 2) the number of features is very large, 3) countermeasures are implemented, and 4) we cannot make assumptions about data distribution. In the beginning, the side-channel community showed most interest in techniques like random forest [23, 24] and support vector machines [25, 26]. More recently, multilayer perceptron [27, 28] and convolutional neural networks [5, 7, 8] emerged as the most potent approaches. Convolutional neural networks were demonstrated to be capable to cope with random delay countermeasure due to their spatial invariance property [7, 8]. At the same time, the fully-connected layers in multilayer perceptron and convolutional neural networks are effective against masking countermeasure as they produce the effect of a higher-order attack (combining features) [8, 29]. Finally, as far as we are aware, the only application of autoencoders for profiled SCA is done by Maghrebi et al., but there the authors use it for classification and not noise removal and they report a relatively poor performance when compared to CNNs [5].

## 1.2 Our Contributions

In this paper, we consider denoising the protected traces with convolutional autoencoder (CAE), which to the best of our knowledge, has not been explored before in the side-channel domain. More precisely, we introduce a novel approach to remove the countermeasures and we provide:

1. A convolutional autoencoder architecture, which requires a limited number of traces to train and can denoise/remove the effect of various hiding techniques.
2. A methodology to recover the ground truth of the traces.
3. An approach to regulate the traces from different sources.

We emphasize that from an attacker perspective, a CAE can be easily trained by noisy(protected)–clean(unprotected) traces pairs. Once the training finishes,

the autoencoder can be used to denoise the leakages from real-world devices. The method proposed in this paper can significantly decrease the protection level of devices that in turn become more vulnerable to the existing attack methods.

This paper is organized as follows. In Section 2, we provide details about profiled SCAs, machine learning techniques we use, and the dataset we investigate. In Section 3, we provide details about denoising autoencoders and their convolutional version. Section 4 gives experimental results when considering the effects of noise either separately or combined. Finally, in Section 5, we conclude the paper and present possible future research directions.

## 2 Background

In this section, we start by introducing the notation we follow. Afterward, we discuss profiled side-channel analysis and neural networks. Finally, we give details about the ASCAD dataset we use in the rest of the paper.

### 2.1 Notation

Let $k^*$ denote the fixed secret cryptographic key (byte), $k$ any possible key hypothesis, and $p$ plaintext. To guess the secret key, the attacker first needs to choose a leakage model $Y(p, k)$ (or $Y$ when there is no ambiguity) depending on the key guess $k$ and some known text $p$, which relates to the deterministic part of the leakage. The size of the keyspace equals $|K|$.

For the autoencoder, we denote its input as $\mathcal{X}$. The encoder of an autoencoder is denoted as $\phi$ and the decoder as $\psi$. Its latent space is denoted as $\mathcal{F}$. As for the training data, we refer to protected leakages/traces (with noise and countermeasures) as noisy leakages/traces; the unprotected leakages are denoted as clean leakages/traces.

### 2.2 Profiled Side-channel Analysis

In the context of implementation attacks, attacks target physical leakage from the insecure implementation of otherwise theoretically secure cryptographic algorithms. The profiled side-channel attacks represent the most powerful category of SCAs as we assume an attacker with access to an open (keys can be chosen/or are known by the attacker) clone device. Then, the attacker can use that clone device to obtain measurements from it and construct a characterization model of the device's behavior. To launch an attack, the attacker then collects a few power traces from the attack device where the secret key is not known. By comparing the attack traces with the characterized model, the secret key can be revealed. Ideally, the secret key can be obtained with a single trace from the attack device. This is difficult in practice due to the effect of the noise, countermeasures, and a finite number of traces in the profiling phase (while we assume the attacker is

not bounded in his power and he can collect any number of traces, that number represent a small fraction of all possible measurements).

To assess the performance of the attacker, one uses a metric denoting the number of measurements required to obtain the secret key. A common example of such a metric is guessing entropy (GE) [30]. GE represents the average number of key candidates an adversary needs to test to reveal the secret key after conducting a side-channel analysis. In particular, given $T$ amount of samples in the attacking phase, an attack outputs a key guessing vector $g = [g1, g2, ..., g_{|K|}]$ in decreasing order of probability. Then, guessing entropy is the average position of $k^*$ in $g$ over several experiments.

## 2.3 Neural Networks

A neural network is an interconnected assembly of simple processing elements, units or nodes, whose functionality is based on the biological process occurring in the brain [31]. In general, a neural network consists of three blocks: an input layer, one or more hidden layers, and an output layer, whose processing ability is represented by the strength (weight) of the inter-unit connections, learning from a set of training patterns from the input layer.

In the supervised machine learning paradigm, neural networks work in two phases: training and testing. In the training phrases, the goal is to learn a function $f$, s.t. $f : \mathcal{X} \to \mathcal{Y}$, given a training set of $N$ pairs $(x_i, y_i)$. Here, for each example (trace) $x$, there is a corresponding label $y$, where $y \in \mathcal{Y}$. Once the function $f$ is obtained, the testing phase starts with the goal to predict the labels for new, previously unseen examples.

**Convolutional Neural Networks** Convolutional neural networks (CNNs) are a type of neural network originally designed for 2-dimensional convolutions as inspired by the biological processes of animals' visual cortex [32]. They are commonly used for image classification but in recent years, they have shown their strengths for time series data [33, 34], speech [35], but also security applications [36].

CNNs resemble ordinary neural networks (e.g., multilayer perceptron) from the architecture perspective: they consist of several layers where each layer is made of neurons. CNN usually consists of three types of layers: convolutional layers, pooling layers, and fully-connected layers. Each layer of a network transforms one volume of activation functions to another through a differentiable function. When considering the CNN architecture, input holds the raw features. Convolution layer computes the output of neurons that are 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. The fully-connected layer computes either the hidden activations or the class scores. To avoid the overfitting, batch normalization layer, which normalizes the input layer by adjusting and scaling the activations is commonly added to the network.

**Autoencoders** Autoencoders were first introduced in the 1980s by Hinton and the PDP group [37] to address the problem of "backpropagation without a teacher". Unlike other neural network architectures that map the relationship between input and its labels, an autoencoder aims to transform inputs into outputs with the least possible amount of distortion [38]. Benefiting from its unsupervised learning characteristic, autoencoder has been used in many applications such as data compression [39], anomaly detection [40], and image recovery [21]. The basic structure of an MLP-based autoencoder is shown in Figure 1.
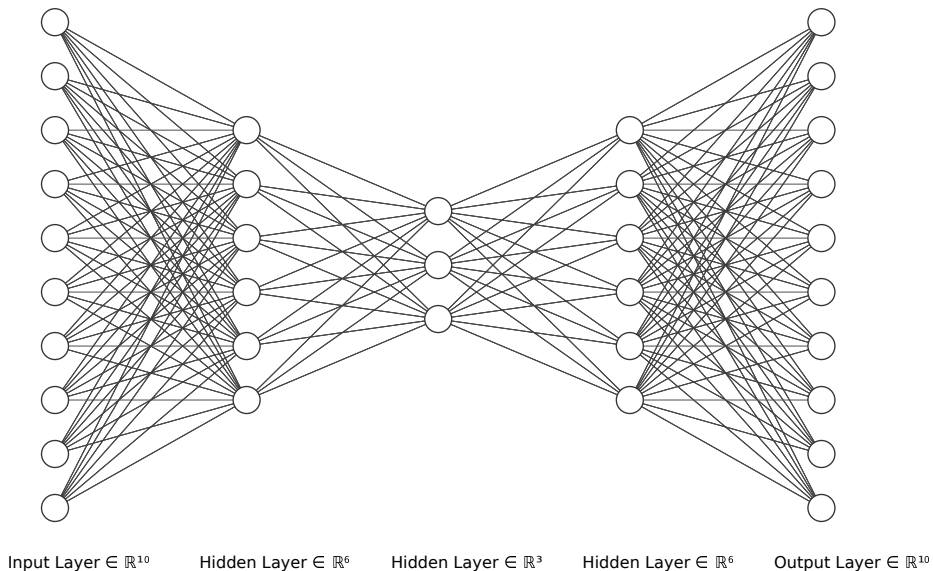


Input Layer $\in \mathbb{R}^{10}$  Hidden Layer $\in \mathbb{R}^{6}$  Hidden Layer $\in \mathbb{R}^{3}$  Hidden Layer $\in \mathbb{R}^{6}$  Output Layer $\in \mathbb{R}^{10}$

Fig. 1: An example of an autoencoder network with three hidden layers (created with NN-SVG [41]).

An autoencoder consists of two parts: encoder ($\phi$) and decoder ($\psi$). Intuitively, the encoder squeezes the input with more features to its bottleneck with fewer features, while the goal of the decoder is to reverse this process. More precisely, the goal of the encoder is to transfer the input to its latent space $\mathcal{F}$, in other words $\phi : \mathcal{X} \to \mathcal{F}$. The decoder, on the other hand, reconstructs the input from latent space, which is equivalent to $\psi : \mathcal{F} \to \mathcal{X}$. To train an autoencoder, the goal is to minimize the distortion when transferring the input to the output (Eq. (1)). Stated differently, the most representative input features are forced to be kept in the smallest layer in the network.

$$\phi, \psi = \arg\min_{\phi, \psi} \|X - (\psi \circ \phi)X\|^2. \tag{1}$$

When applying the autoencoder for the denoising purpose, the input and output are not identical but represent the noisy-clean data pairs. Ideally, a well-

trained autoencoder can keep the most representative information (e.g., characteristics of the original data) in its bottleneck while neglecting the less important features (e.g., random noise). Consequently, the original data (without noise) can be recovered by feeding noisy data to the input of the autoencoder.

A similar idea can also be applied to remove the countermeasures from the leakage traces. Since the noise and countermeasures existing in the leakage are random in the amplitude or time domain, we can also consider them as noise. For example, we can refer to the white noise as the noise in the amplitude domain; the global jitters (desynchronization), random delay interrupts, and clock jitters, on the other hand, can be considered as the noise in the time domain. Since the autoencoder is good at extracting the important features while neglecting randomness, we use autoencoder to filter out this noise and recover the ground truth of the traces. Then, one can expect that with the recovered traces the attack efficiency will be dramatically improved.

## 2.4  The ASCAD dataset

The ASCAD dataset was introduced by Prouff et al. to provide a benchmark to evaluate machine learning techniques in the context of side-channel attacks [29]. An ATMega8515 device was used to record the emitted EM radiation during the execution of a software AES implementation protected by known masks. All traces were captured with a sensor attached to an oscilloscope sampling at 2 $GS/s$.

There are two data sets recorded in different conditions: fixed key encryption and random key encryption. For the data with fixed key encryption, the dataset provided separate HDF5 files with different synchronization level: ASCAD.h5, ASCAD_desync50.h5, and ASCAD_desync100.h5. The traces in the ASCAD.h5 file are time-aligned in a prepossessing step, whereas the traces in ASCAD_desync50.h5 and ASCAD_desync100.h5 have been shifted with a maximum jitters window of respectively 50 and 100 samples [29]. Each file contains 60 000 EM traces (50 000 training / cross-validation traces and 10 000 test traces). Each trace consists of 700 points of interest.

For the data with random key encryption, there are 200 000 traces in the profiling dataset that is provided to train the (deep) neural network models. A 100 000 traces attack dataset is used to check the performance of the trained models after the profiling phase. A window of 1 400 points of interest is extracted around the leaking spot.

Throughout the paper, we use the raw traces and the pre-selected window of relevant samples per trace corresponding to masked S-box for $i = 3$. As a leakage model, we use the unprotected S-box output, i.e.:

$$Y(i) = \texttt{Sbox}[(p[i] \oplus k[i])]. \tag{2}$$

Note that the model given in Eq. (2) does not leak information directly as it is first-order protected. Consequently, we do not state a model-based SNR. The SNR for the ASCAD dataset is $\approx 0.8$ under the assumption we know the mask (shown in Figure 4a) while it is almost 0 with the unknown mask.

# 3    Denoising with Convolutional Autoencoder

In this section, we discuss the attacker model and afterward, the details about the convolutional autoencoder architecture.

## 3.1    Denoising Strategy

As discussed in Section 2.3, noisy-clean trace pairs are required to train a denoising autoencoder. Inspired by the profiled side-channel analysis, we devise a denoising strategy shown in Figure 2.
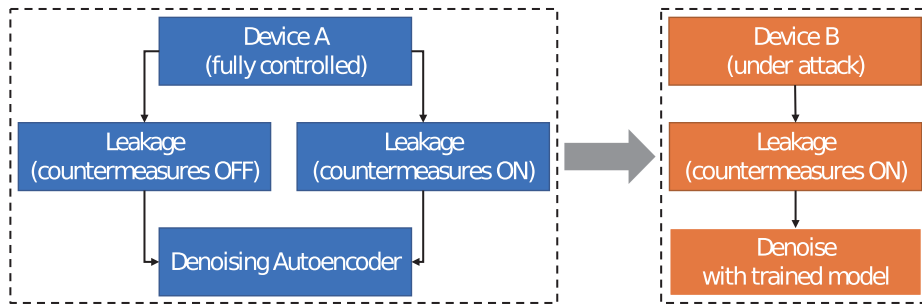


Fig. 2: Denoising strategy.

We assume an attacker has full control of a device (Device A). Specifically, he can enable/disable the implemented countermeasure. To attack the real devices with countermeasures enabled (Device B), he first acquires several traces with and without countermeasures from Device A to build the training sets. Then the attacker uses these traces to train the denoising autoencoder. Once the training process is finished, the trained model can be used to pre-process the leakage traces obtained from Device B. Finally, with "cleaner" traces reconstructed by the denoising autoencoder, an attacker can eventually retrieve the secret information with less effort.

## 3.2    Convolutional Autoencoder Architecture

An autoencoder can be implemented with different neural network architectures, the most common examples are the MLP-based autoencoder and convolutional autoencoder (CAE). Since related works indicate that CNN outperforms MLP in dealing with the side-channel leakages [5,28], we use the convolution layer as the basic element for denoising purpose.

To maximize the denoising ability of the proposed architecture, we tune the hyperparameters by evaluating the CAE performance towards different types of noise. The tuning range and selected hyperparameters are shown in Table 1.

We use the *SeLU* activation function to avoid vanishing and exploding gradient problems [42] and *He Uniform* initialization to improve weight initialization [43].

Table 1: CAE hyperparameter tuning.

| Hyperparameter | Range | Selected |
|---|---|---|
| Optimizer | Adam, RMSProb, SGD | RMSProb |
| Weight initialization | Uniform distribution, He uniform | He uniform |
| Activation function | tanh, ReLU, SeLU | SeLU |
| Learning Rate | 1e-5, 5e-5, 1e-4, 5e-4 | 1e-4 |
| Batch size | 32, 64, 128, 256 | 128 |
| Epochs | 30, 50, 70, 100, 200 | 100 |
| Training sets | 1 000, 5 000, 10 000, 20 000 | 10 000 |
| Validation sets | 1 000, 2 000 | 2 000 |

In terms of the autoencoder architecture, we observed that when dealing with trace desynchronization, an autoencoder that has shallow architecture can denoise the traces successfully. Still, when we introduce other types of noise into the traces while keeping the same hyperparameters, such autoencoder cannot recover the ground truth of the traces. Therefore, we decided to increase the depth of the autoencoder to ensure it will be suitable for different types of noise.

The size of the latent representation in the middle of the autoencoder is a critical parameter that should be fine-tuned. One should be aware that although the autoencoder can reconstruct the input, some information from the input is lost. For the denoising purpose, we aim at maximizing the removal of noise while minimizing the loss of useful information. By choosing a smaller size of the bottleneck, the signal quality will be degraded. In contrast, a larger size of bottleneck may introduce less critical features to the output. To better control the size of the latent space, we flatten the output of the convolutional blocks and introduce a fully-connected layer as the middle layer in our proposed architecture.

The CAE architecture used for this paper is shown in Table 2. The convolution block (*Convblock* in the table) normally consists of three layers: convolution layer, activation layer (function), and Max pooling layer. As we noticed that an autoencoder implemented in this manner suffers from overfitting and poor performance in denoising the validation traces, we add the batch normalization layer in each convolution block.

The convolutional encoder converts the leakages to its latent representation by passing through multiple convolution blocks. Then, these compressed high-dimensional features are further calculated with a fully-connected layer, which is then reshaped to the convolution size and reconstructed to the target traces. Note that the latent space size is controlled by the number of neurons in the fully-connected layer. To ensure the CAE output has the same shape with the clean

traces, the size of the fully-connected layer $S_{latent}$ is calculated with Eq. (3):

$$S_{latent} = \frac{S_{clean}}{\prod_{i=1}^{n} S_{pooli}} * N_{filter0}. \tag{3}$$

$S_{clean}$ is the size of the target clean traces, $S_{pooli}$ represents the $i$th non-zero pooling stride of the decoder, and $N_{filter0}$ represents the number of the filters of the first Deconv layer.

Table 2: CAE architecture.

| Block/Layer | Filter size | Filter number | Pooling stride |
|---|---|---|---|
| Conv block * 2 | 2 | 256 | 0 |
| Conv block | 2 | 256 | 5 |
| Conv block * 2 | 2 | 128 | 0 |
| Conv block | 2 | 128 | 2 |
| Conv block * 2 | 2 | 64 | 0 |
| Conv block | 2 | 64 | 2 |
| Flatten | - | - | - |
| Fully-connected | - | - | - |
| Reshape | - | - | - |
| Deconv block | 2 | 64 | 2 |
| Deconv block * 2 | 2 | 64 | 0 |
| Deconv block | 2 | 128 | 2 |
| Deconv block * 2 | 2 | 128 | 0 |
| Deconv block | 2 | 256 | 5 |
| Deconv block * 2 | 2 | 256 | 0 |
| Deconv block | 2 | 1 | 0 |

## 4   Experimental Results

Several noise types can exist in the power profiles or EM traces. To show the effectiveness of the proposed CAE in dealing with different types of noise, we investigate four types of noise/countermeasures that commonly exist in the devices: Gaussian noise, desynchronization (misalignment), random delay interrupts (RDI), and clock jitters. All of these noise/countermeasures are simulated based on the observation or implementation of the real devices. To compare the denoising performance of CAE with the existing techniques commonly used by attackers, the denoising performance of trace averaging is also evaluated in this paper. Here, 10 traces are averaged to obtain a single trace.

Throughout the experiments, we use the ASCAD dataset (with fixed key and random key) as the training datasets. Note that the ASCAD dataset is

masked and there is no first-order leakage. To obtain a direct intuition of the properties of the reconstructed traces, we evaluate the SNR of the masked S-box output (Eq. (4)). The SNR is sometimes named F-Test to refer to its original introduction in [44]. For a noisy observation $L_t$ at each time sample $t$ of an event $Z$, it is defined in Eq. (5):

$$Y(i) = \texttt{Sbox}[(p[i] \oplus k[i]) \oplus r_{out}]. \tag{4}$$

$$\text{SNR}\,[t] \triangleq \frac{\text{Var}\,\mathbb{E}[L_t|Z]}{\mathbb{E}\,\text{Var}[L_t|Z]}. \tag{5}$$

The numerator denotes the signal magnitude and the denominator denotes the noise magnitude estimate. The model for the attacking is the CNN_best as given in the ASCAD paper [29]. The selected hyperparameters are shown in Table 3. The quality of the recovered traces is evaluated by guessing entropy (GE). For a good estimation of GE, the attack traces are randomly shuffled and 100 GEs are computed to obtain the average value.

Table 3: CNN Hyperparameter.

| Hyperparameter | Value |
| --- | --- |
| Optimizer | RMSProb |
| Weight initialization | Uniform distribution |
| Activation function | ReLU |
| Learning Rate | 1e-5 |
| Batch size | 200 |
| Epochs | 100 |
| Training sets | 35 000 |
| Validation sets | 5 000 |

### 4.1  Baseline

We first attack the synchronized ASCAD dataset (ASCAD.h5) as the baseline in the comparison with different scenarios. The detailed attack approach can be found in the original paper [29]. An example of two traces to be attacked is presented in Figure 3. From the comparison of those two traces, we see that most of the parts are identical, but there are still some variations existing (especially on the peaks).

As mentioned, we use the CNN_best model for the attack. The SNR of the masked S-box output and the resulting guessing entropy are presented in Figure 4. Without noise and countermeasures in the traces, guessing entropy reaches zero after 200 traces, indicating that the real key has been successfully obtained. The outcomes of this attack are used as the baseline and the following attacks on noisy and denoised traces are evaluated based on these results.
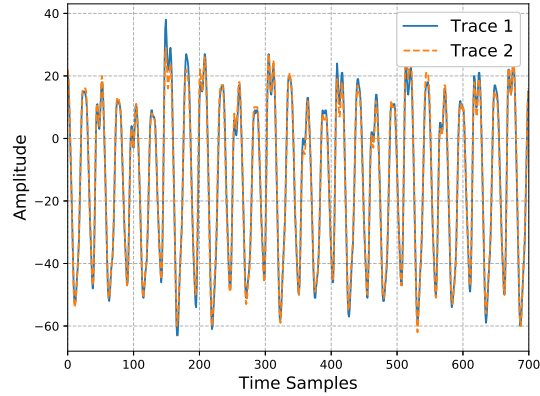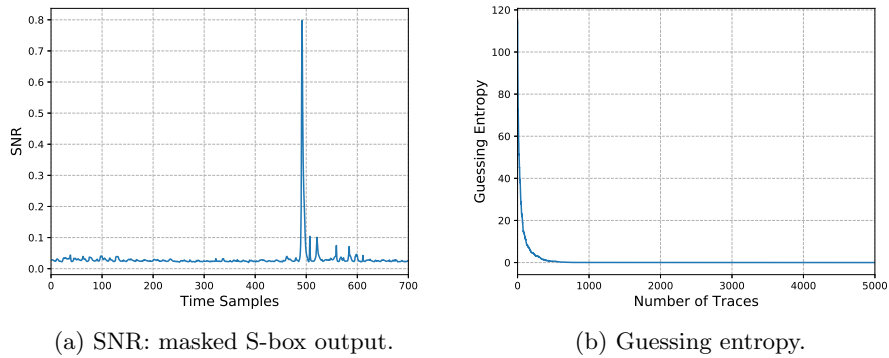
Fig. 3: Baseline: example traces.



(a) SNR: masked S-box output.

(b) Guessing entropy.

Fig. 4: Baseline: SNR and guessing entropy.

## 4.2 Gaussian Noise

The Gaussian noise is the most common type of noise existing in side-channel traces. The transistor, data buses, the transmission line to the record devices such as oscilloscopes or even the work environment can be the source of Gaussian noise. The noise can also be artificially introduced by dummy operation or dedicated noise engine. In terms of trace leakage, the increment of the noise level hides the useful patterns and reduces the signal-to-noise (SNR) ratio. Consequently, the noise influences the effectiveness of an attack, i.e., more traces are needed to obtain the attacked intermediate data.

To demonstrate the influence of the Gaussian noise, we simulate it by adding a random number to each point of the trace. Specifically, a random number is uniformly distributed between -20 to 20 to simulate the real leakage behavior with the Gaussian noise. An example of the manipulated trace and its zoom-in

view is shown in Figure 5. Compared with the baseline traces, the Gaussian noise significantly distorted the shape of the original traces in the amplitude domain.



(a) Gaussian noise: example traces.

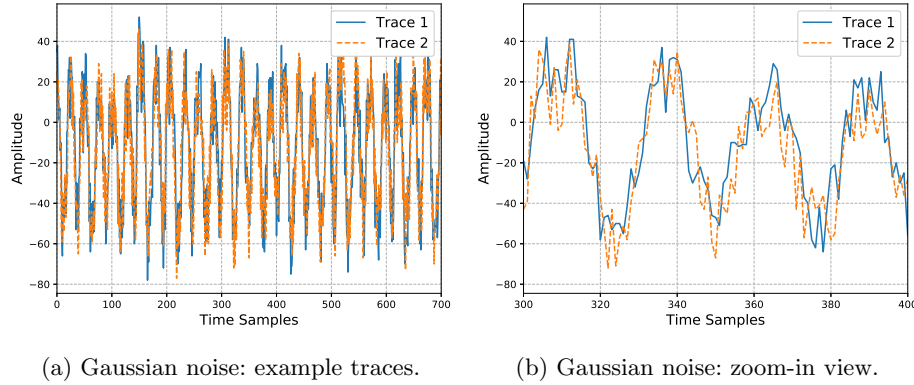(b) Gaussian noise: zoom-in view.

Fig. 5: Gaussian noise: example traces and its zoom-in view.

Then, we denoise the Gaussian noise with trace averaging as well as CAE proposed in this paper. The SNR of the noisy and denoised traces are shown in Figure 6a while GE values after deep learning attacks are shown in Figure 6b.
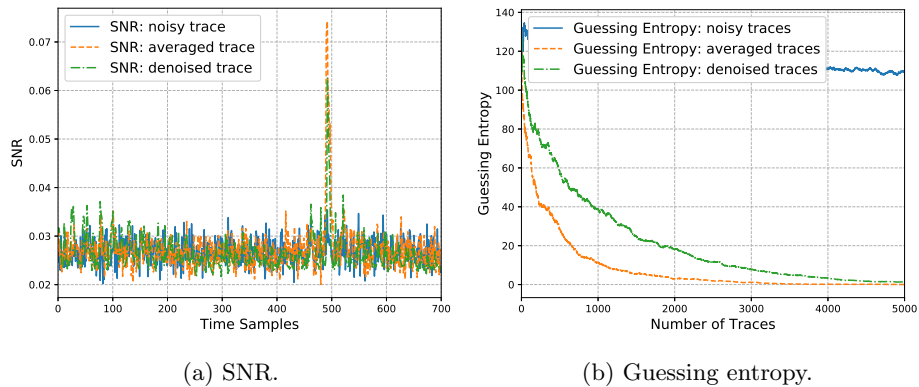


(a) SNR.

(b) Guessing entropy.

Fig. 6: Noisy and denoised traces (Gaussian noise): SNR and guessing entropy.

Firstly, when considering the baseline attack with "clean" traces, the existence of the Gaussian noise reduces the leakage of the attacked intermediate data: the SNR drops from 0.8 to 0.03. Fortunately, the SNR of denoised traces clearly shows that the CAE has filtered out a part of the Gaussian noise, as the

SNR peak of the denoised traces is more than three times higher than the noisy one.

From the attack (GE) perspective, GE converges in both cases when the number of trace increases. For the noisy traces, 5 000 traces reduces GE from 130 to only 110. On the other hand, outstanding attacking performance is obtained with denoised traces: 5 000 traces are sufficient to reach GE of 0 when denoised with both CAE and averaging techniques. It is worth to note that GE of averaged traces is slightly better than GE of CAE, proving that trace averaging is a good candidate in removing the Gaussian noise. Still, we can conclude that CAE can remove the Gaussian noise and consequently improve the attacking efficiency.

### 4.3   Desynchronization

Well-synchronized traces can significantly improve the correlation of the intermediate data. The alignment of the traces is, therefore, an essential step for the side-channel attack. To align the traces, normally, an attacker should select a distinguishable trigger/pattern from the traces, so that the following part can be aligned using the selected part as a reference. There are two limitations to this approach. First, the selected trigger/pattern should be distinctive, so that it will not be obfuscated with other patterns and lead to misalignment. Second, due to the existence of the signal jitters and other countermeasures, the selected trigger should be sufficiently close to the points of interest, thus minimizing the noise effect. From a practical point of view, a good reference that meets both limitations is not always easy to find. Even with an unprotected device, sometimes the traces synchronization can be a challenging task.

We consider the desynchronization as a type of noise existing in the traces. Different from the Gaussian noise, the desynchronization noise adds randomness to the time domain. To show the effect of the traces desynchronization, we use traces with a maximum of 50 points of desynchronization (ASCAD_desync50 [29]). An example of the traces with desynchronization is shown in Figure 7.

Next, we attack the misaligned traces as well as the denoised traces from CAE. The SNR and GE results are shown in Figure 8.

Figure 8a shows that only CAE can increase the SNR of the intermediate data (0.16). The averaged traces, on the other hand, have similar SNR (0.03) with noisy ones, indicating that traces averaging is ineffective in dealing with desynchronization.

From Figure 8b, GE of the noisy traces converges faster than the averaged traces. From this result, first, we conclude that the averaging of the desynchronized traces is harmful to recovering the traces. Second, CNN proves its ability to delimit the desynchronization effect, as GE converges to around 55 with 5 000 traces when attacking the noisy traces. Still, considering that the original "clean" traces only needed around 200 traces to retrieve the key, the desynchronization indeed degraded the performance of the attack. One can expect that performance to become even worse with an increased desynchronization level.

CAE provides an alternative simple approach in synchronizing the traces. By training a CAE with desynchronized–synchronized traces pair, the model
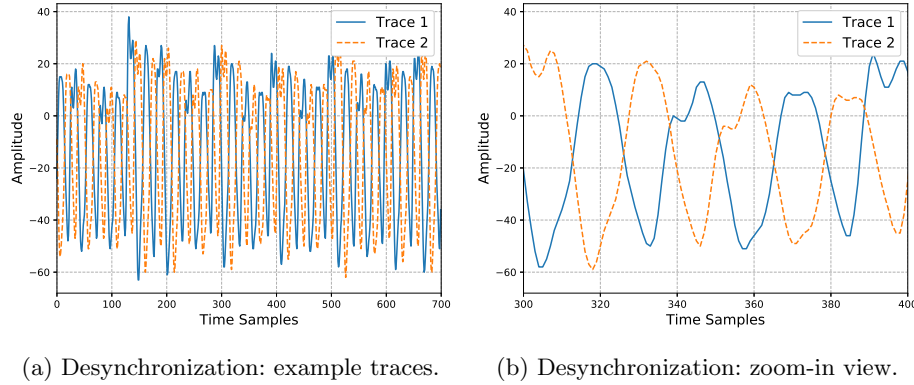
(a) Desynchronization: example traces.

(b) Desynchronization: zoom-in view.

Fig. 7: Desynchronization: example traces and its zoom-in view.
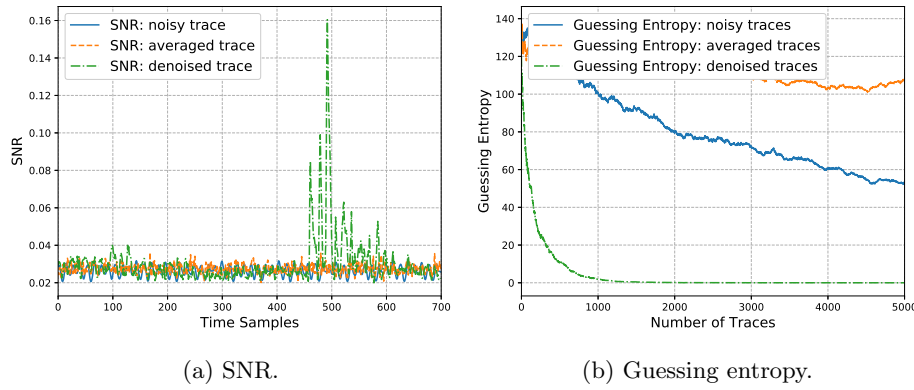


(a) SNR.

(b) Guessing entropy.

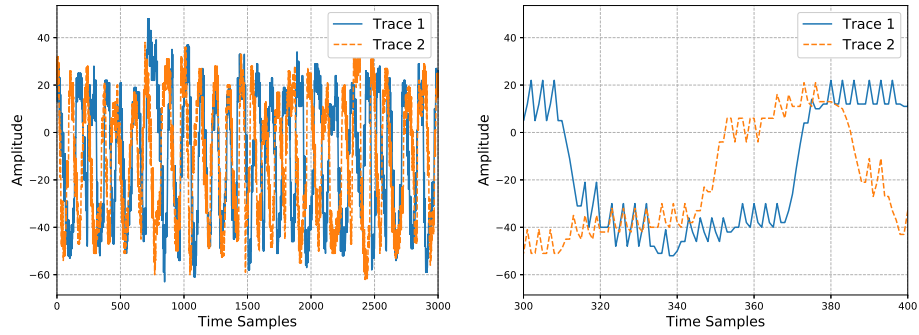Fig. 8: Noisy and denoised traces (desynchronization): SNR and guessing entropy.

can align the traces automatically. Consequently, the number of required traces to retrieve the key reduces to around 1 000, which performs the best from all three considered scenarios.

### 4.4 Random Delay Interrupts (RDIs)

Desynchronization introduces the global time-randomness to the entire trace. RDIs, on the other hand, lead to the time-randomness locally. As a type of countermeasure normally implemented in the software, the existence of RDIs breaks the traces into fragments, thus significantly increasing the randomness of traces in the time domain and reducing the correlation of the attacked intermediate data.

To simulate the RDIs, we inject RDIs based on the Floating Mean method introduced in [45]. The RDIs implemented in such a way can provide more variance to the traces when compared with the uniform RDI distribution, thus further
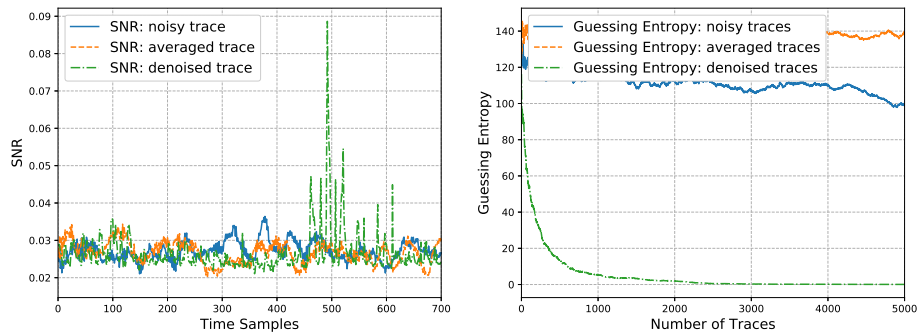
increasing the attack difficulty. An example of the traces with RDIs with its zoom-in view is shown in Figure 9. We observe that more randomness was introduced locally to the traces when compared to the traces with desynchronization.



(a) Random Delay Interrupts: example traces.

(b) Random Delay Interrupts: zoom-in view.

Fig. 9: Random Delay Interrupts: example traces and its zoom-in view.

As a result, the SNR and rank of the original traces with RDIs and denoised traces with averaging and CAE are shown in Figure 10. The guessing entropy of the traces with RDIs converges slowly (120 to 100 with 5 000 traces), indicating that the CNN_best model we are using is not powerful enough to extract the useful patterns and retrieve the key with 5 000 traces. We can conclude that RDIs implemented in this way dramatically increase the attack difficulty.



(a) SNR.

(b) Guessing entropy.

Fig. 10: Noisy and denoised traces (RDIs): SNR and guessing entropy.

Trace averaging is again not helpful to remove the countermeasure and the GE value fluctuates around 140. By applying the CAE, the effect of RDIs has been reduced dramatically: the SNR increases three times to 0.09 and GE reaches zero after around 2 500 traces. From the results, we can conclude that CAE can recover the original traces from the noisy traces with RDIs countermeasure.

## 4.5 Clock Jitters

Clock jitters is a classical hardware countermeasure against side-channel attacks, realized by introducing the instability in the clock [7]. Comparable to the Gaussian noise that introduces randomness to every point in the amplitude domain, the clock jitters increase the randomness for each point in the time domain. Indeed, the accumulation of the deforming effect increase the misalignment of the traces as well as decrease the correlation of the intermediate data. As a consequence, the attacked intermediate data become more difficult to retrieve. Here, we simulate the clock jitters by randomly adding or removing points with a predefined range. More precisely, we generate a random number $r$ that is uniformly distributed between -4 to 4 to simulate the clock variation in a magnitude of 16. When scanning each point in the trace, $r$ points will be added to the trace if $r$ is greater than zero. Otherwise, the following $r$ points in the trace are deleted. An example of the traces with clock jitters is shown in Figure 11.



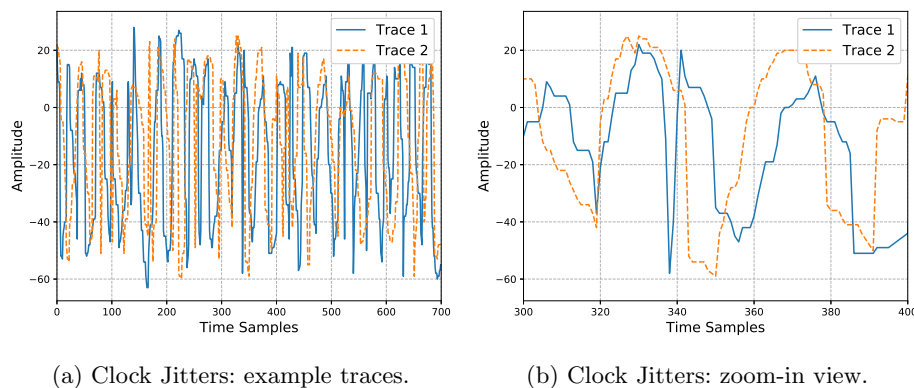(a) Clock Jitters: example traces.  (b) Clock Jitters: zoom-in view.

Fig. 11: Clock Jitters: example traces and its zoom-in view.

We denoise the clock jitters with CAE and a comparison of the attack results for the noisy and denoised traces with averaging and CAE is shown in Figure 12. As in the previous settings, traces averaging is not efficient in dealing with time-based noise and countermeasures. The proposed CAE, on the other hand, successfully reduces the effect of clock jitters. From the SNR perspective, although the denoised SNR is the lowest (0.044) for all four different types of noise, it still outperforms its noisy and averaged counterparts. When compared

with noisy traces, the GE of the denoised traces converges to around 20 after applying 5 000 traces.



(a) SNR.                                    (b) Guessing entropy.
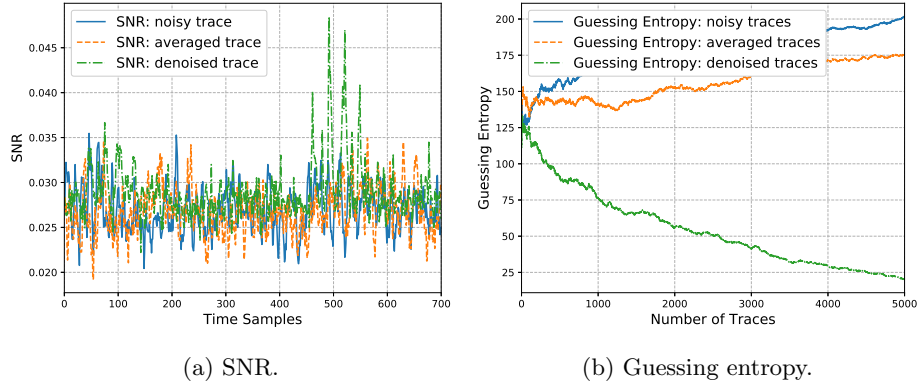
Fig. 12: Noisy and denoised traces (clock jitters): SNR and guessing entropy.

To conclude, the proposed CAE proves its ability to delimit the effect of the Gaussian noise, desynchronization, random delay interrupts, and clock jitters. Traces averaging performs well in removing the Gaussian noise but is ineffective in dealing with the noise and countermeasures in the time domain.

### 4.6   Combining the Effects of Gaussian Noise and Countermeasures

In the previous section, we add and denoise different types of noise individually. Now, we investigate a more realistic situation by adding all of the discussed noise types together and then verifying the effectiveness of the CAE approach. Here, we test two different datasets: AES with a fixed key and AES with random keys. Since trace alignment is proved to be inefficient in dealing with time-based noise and countermeasure, we only evaluate the SNR and GE of the noisy and denoised traces with CAE.

**AES with Fixed Key**  Similar to the procedure of the previous sections, we calculated the SNR and GE of the noisy and denoised traces and make a comparison between them. The SNR comparison is presented in Figure 13a. Compared with the noisy traces, the SNR of denoised traces is increased sightly (0.0345). This observation indicates that the combination of the noise types degraded the performance of the CAE.

From the GE plot (Figure 13b), the noisy traces do not converge with the increasing number of traces. The GE of denoised traces, on the other hand, reaches around 30 with 5 000 traces. Although it converges slower than the denoised traces with a single type of noise, CAE still proves its capability in removing the combined effect of noise and countermeasures.
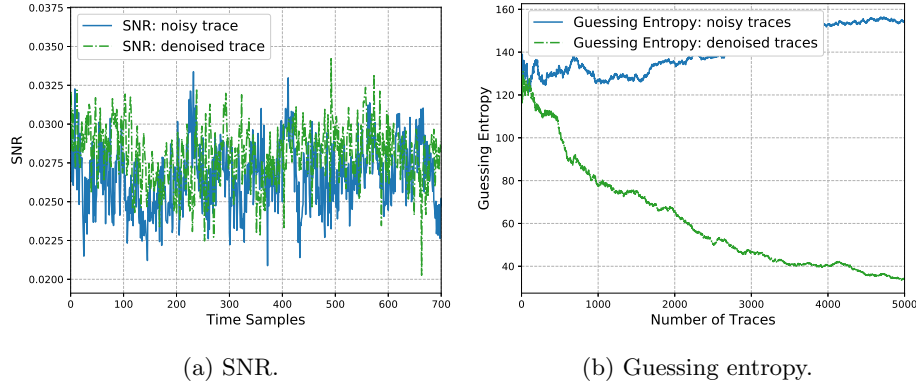
(a) SNR.

(b) Guessing entropy.

Fig. 13: Noisy and denoised traces (all): SNR and guessing entropy.

**AES with Random Keys** Next, we verify the performance of the CAE by trying to denoise the AES traces with random keys. To retrieve the correct key from the leakage traces, we first train the model with leakage with a random but known key, then use the trained model to attack the leakages and try to retrieve the unknown key. When comparing with the fixed-key traces, the randomness of the key introduces more variance into the traces, thus further increasing the difficulties in denoising the traces. In terms of attack settings, we use 1 400 POI for the attack. The attacked intermediate data was kept the same.
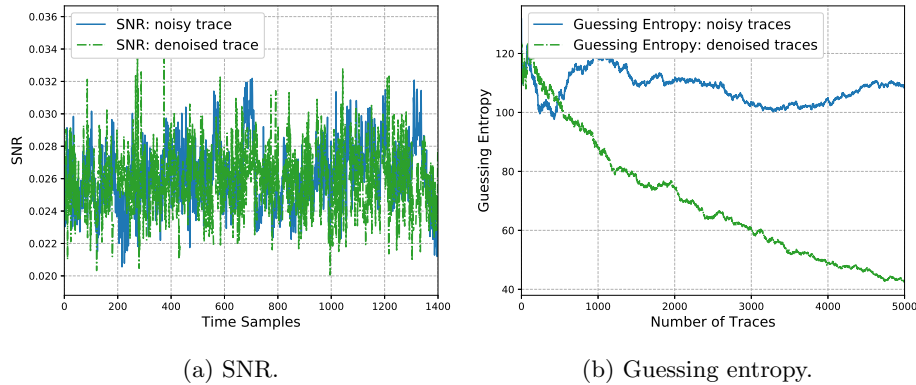


(a) SNR.

(b) Guessing entropy.

Fig. 14: Noisy and denoised traces with variable key (all): SNR and guessing entropy.

Based on the results, there is no significant improvement from the SNR perspective. Still, guessing entropy indicates the improved performance as a result of CAE: GE value converged to around 43 with 5 000 traces. The GE of the

noisy traces, on the other hand, fluctuates above 100 regardless of the number of traces. Therefore, we can conclude that the proposed CAE can denoise the leakage in both fixed key and variable key scenarios.

## 5   Conclusions and Future Work

In this paper, we introduce a convolutional autoencoder to remove the noise and countermeasure from the leakage traces. Different types of noise and counter-measures, such as Gaussian noise, desynchronization, random delay interrupts, and clock jitters are simulated and attacked with the CNN model. The noisy traces and denoised traces are compared from two different perspectives: SNR and guessing entropy. We show that the CAE can reconstruct the ground truth from the noisy traces and thus improve the attack efficiency. Additionally, we simulate the scenario where all noise types and countermeasures are combined into the measurements. Surprisingly, the proposed CAE can still remove the noise and find out the underlying ground truth. Throughout the paper, two types of leakage traces (one encrypted with fixed and another with random keys) are tested and in both cases, the CAE proves its performance for denoising purpose.

CAE provides an attacker with a powerful tool to pre-process the traces. Besides that, we expect the denoising autoencoder could be used to solve other problems like portability [17]. For the attack portability, the biggest obstacle comes from the variance among different devices. These variances introduce the variation of the trace, making the attack model generated from one device difficult to transfer to another one. With the help of the autoencoder, this problem can be solved in another way: instead of working with the training model, we can consider the traces variation as noise and use denoising autoencoder to remove it. The basic idea is to select the traces from one device as a reference and then train an autoencoder to convert the traces from other devices to the selected device. Since the variation of the traces between devices should be diminished with the help of the autoencoder, the attack model from one device could be applied to other devices.

## References

1. Mangard, S., Oswald, E., Popp, T.: Power Analysis Attacks: Revealing the Secrets of Smart Cards. Springer (December 2006) ISBN 0-387-30857-1, http://www.dpabook.org/.
2. Quisquater, J.J., Samyde, D.: Electromagnetic analysis (ema): Measures and counter-measures for smart cards. In Attali, I., Jensen, T., eds.: Smart Card Programming and Security, Berlin, Heidelberg, Springer Berlin Heidelberg (2001) 200–210
3. Kocher, P.C., Jaffe, J., Jun, B.: Differential power analysis. In: Proceedings of the 19th Annual International Cryptology Conference on Advances in Cryptology. CRYPTO '99, London, UK, UK, Springer-Verlag (1999) 388–397
4. Chari, S., Rao, J.R., Rohatgi, P.: Template Attacks. In: CHES. Volume 2523 of LNCS., Springer (August 2002) 13–28 San Francisco Bay (Redwood City), USA.

5. Maghrebi, H., Portigliatti, T., Prouff, E.: Breaking cryptographic implementations using deep learning techniques. In: International Conference on Security, Privacy, and Applied Cryptography Engineering, Springer (2016) 3–26

6. Picek, S., Samiotis, I.P., Kim, J., Heuser, A., Bhasin, S., Legay, A.: On the performance of convolutional neural networks for side-channel analysis. In Chattopadhyay, A., Rebeiro, C., Yarom, Y., eds.: Security, Privacy, and Applied Cryptography Engineering, Cham, Springer International Publishing (2018) 157–176

7. Cagli, E., Dumas, C., Prouff, E.: Convolutional neural networks with data augmentation against jitter-based countermeasures. In: International Conference on Cryptographic Hardware and Embedded Systems, Springer (2017) 45–68

8. Kim, J., Picek, S., Heuser, A., Bhasin, S., Hanjalic, A.: Make some noise. unleashing the power of convolutional neural networks for profiled side-channel analysis. IACR Transactions on Cryptographic Hardware and Embedded Systems (2019) 148–179

9. Brier, E., Clavier, C., Olivier, F.: Correlation power analysis with a leakage model. In: International Workshop on Cryptographic Hardware and Embedded Systems, Springer (2004) 16–29

10. Picek, S., Heuser, A., Jovic, A., Batina, L.: A systematic evaluation of profiling through focused feature selection. IEEE Transactions on Very Large Scale Integration (VLSI) Systems **27**(12) (Dec 2019) 2802–2815

11. Thiebeauld, H., Gagnerot, G., Wurcker, A., Clavier, C.: Scatter: A new dimension in side-channel. In: International Workshop on Constructive Side-Channel Analysis and Secure Design, Springer (2018) 135–152

12. Chari, S., Jutla, C.S., Rao, J.R., Rohatgi, P.: Towards sound approaches to counteract power-analysis attacks. In: Annual International Cryptology Conference, Springer (1999) 398–412

13. Barthe, G., Dupressoir, F., Faust, S., Grégoire, B., Standaert, F.X., Strub, P.Y.: Parallel implementations of masking schemes and the bounded moment leakage model. In: Annual International Conference on the Theory and Applications of Cryptographic Techniques, Springer (2017) 535–566

14. Clavier, C., Coron, J.S., Dabbous, N.: Differential power analysis in the presence of hardware countermeasures. In: International Workshop on Cryptographic Hardware and Embedded Systems, Springer (2000) 252–263

15. Tiri, K., Verbauwhede, I.: Securing encryption algorithms against dpa at the logic level: Next generation smart card technology. In: International Workshop on Cryptographic Hardware and Embedded Systems, Springer (2003) 125–136

16. Coron, J., Kizhvatov, I.: An Efficient Method for Random Delay Generation in Embedded Software. In: Cryptographic Hardware and Embedded Systems - CHES 2009, 11th International Workshop, Lausanne, Switzerland, September 6-9, 2009, Proceedings. (2009) 156–170

17. Bhasin, S., Chattopadhyay, A., Heuser, A., Jap, D., Picek, S., Shrivastwa, R.R.: Mind the portability: A warriors guide through realistic profiled side-channel analysis. Cryptology ePrint Archive, Report 2019/661 (2019) https://eprint.iacr.org/2019/661.

18. Wei, L., Luo, B., Li, Y., Liu, Y., Xu, Q.: I know what you see: Power side-channel attack on convolutional neural network accelerators. In: Proceedings of the 34th Annual Computer Security Applications Conference, ACM (2018) 393–406

19. Zheng, Y., Zhou, Y., Yu, Z., Hu, C., Zhang, H.: How to compare selections of points of interest for side-channel distinguishers in practice? In Hui, L.C.K., Qing, S.H., Shi, E., Yiu, S.M., eds.: Information and Communications Security, Cham, Springer International Publishing (2015) 200–214

20. Zaid, G., Bossuet, L., Habrard, A., Venelli, A.: Methodology for efficient cnn architectures in profiling attacks. IACR Transactions on Cryptographic Hardware and Embedded Systems **2020**(1) (Nov. 2019) 1–36

21. Gondara, L.: Medical image denoising using convolutional denoising autoencoders. In: 2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW), IEEE (2016) 241–246

22. Lerman, L., Poussier, R., Bontempi, G., Markowitch, O., Standaert, F.X.: Template attacks vs. machine learning revisited (and the curse of dimensionality in side-channel analysis). In: International Workshop on Constructive Side-Channel Analysis and Secure Design, Springer (2015) 20–33

23. Lerman, L., Medeiros, S.F., Bontempi, G., Markowitch, O.: A Machine Learning Approach Against a Masked AES. In: CARDIS. Lecture Notes in Computer Science, Springer (November 2013) Berlin, Germany.

24. Heuser, A., Picek, S., Guilley, S., Mentens, N.: Side-channel analysis of lightweight ciphers: Does lightweight equal easy? In: Radio Frequency Identification and IoT Security - 12th International Workshop, RFIDSec 2016, Hong Kong, China, November 30 - December 2, 2016, Revised Selected Papers. (2016) 91–104

25. Heuser, A., Zohner, M.: Intelligent Machine Homicide - Breaking Cryptographic Devices Using Support Vector Machines. In Schindler, W., Huss, S.A., eds.: COSADE. Volume 7275 of LNCS., Springer (2012) 249–264

26. Picek, S., Heuser, A., Jovic, A., Ludwig, S.A., Guilley, S., Jakobovic, D., Mentens, N.: Side-channel analysis and machine learning: A practical perspective. In: 2017 International Joint Conference on Neural Networks, IJCNN 2017, Anchorage, AK, USA, May 14-19, 2017. (2017) 4095–4102

27. Gilmore, R., Hanley, N., O'Neill, M.: Neural network based attack on a masked implementation of aes. In: 2015 IEEE International Symposium on Hardware Oriented Security and Trust (HOST), IEEE (2015) 106–111

28. Picek, S., Heuser, A., Jovic, A., Bhasin, S., Regazzoni, F.: The curse of class imbalance and conflicting metrics with machine learning for side-channel evaluations. IACR Trans. Cryptogr. Hardw. Embed. Syst. **2019**(1) (2019) 209–237

29. Benadjila, R., Prouff, E., Strullu, R., Cagli, E., Dumas, C.: Study of deep learning techniques for side-channel analysis and introduction to ascad database. ANSSI, France & CEA, LETI, MINATEC Campus, France. Online verfügbar unter https://eprint. iacr. org/2018/053. pdf, zuletzt geprüft am **22** (2018) 2018

30. Standaert, F.X., Malkin, T.G., Yung, M.: A unified framework for the analysis of side-channel key recovery attacks. In: Annual international conference on the theory and applications of cryptographic techniques, Springer (2009) 443–461

31. Gurney, K.: An introduction to neural networks. CRC press (2014)

32. LeCun, Y., Bengio, Y., et al.: Convolutional networks for images, speech, and time series. The handbook of brain theory and neural networks **3361**(10) (1995) 1995

33. Oord, A.v.d., Dieleman, S., Zen, H., Simonyan, K., Vinyals, O., Graves, A., Kalchbrenner, N., Senior, A., Kavukcuoglu, K.: Wavenet: A generative model for raw audio. arXiv preprint arXiv:1609.03499 (2016)

34. Choi, K., Fazekas, G., Sandler, M., Cho, K.: Convolutional recurrent neural networks for music classification. In: 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE (2017) 2392–2396

35. Palaz, D., Collobert, R., et al.: Analysis of cnn-based speech recognition system using raw speech as input. Technical report, Idiap (2015)

36. Sirinam, P., Imani, M., Juarez, M., Wright, M.: Deep fingerprinting: Undermining website fingerprinting defenses with deep learning. In: Proceedings of the

2018 ACM SIGSAC Conference on Computer and Communications Security, ACM (2018) 1928–1943

37. Rumelhart, D.E., Hinton, G.E., Williams, R.J.: Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science (1985)
38. Baldi, P.: Autoencoders, unsupervised learning, and deep architectures. In: Proceedings of ICML workshop on unsupervised and transfer learning. (2012) 37–49
39. Theis, L., Shi, W., Cunningham, A., Huszár, F.: Lossy image compression with compressive autoencoders. arXiv preprint arXiv:1703.00395 (2017)
40. An, J., Cho, S.: Variational autoencoder based anomaly detection using reconstruction probability. Special Lecture on IE **2**(1) (2015)
41. LeNail: NN-SVG: Publication-Ready Neural Network Architecture Schematics. Journal of Open Source Software **4**(33) (2019) 747
42. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. In: Advances in neural information processing systems. (2017) 971–980
43. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE international conference on computer vision. (2015) 1026–1034
44. Fisher, R.A.: On the mathematical foundations of theoretical statistics. Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character **222**(594-604) (1922) 309–368
45. Coron, J.S., Kizhvatov, I.: An efficient method for random delay generation in embedded software. In: International Workshop on Cryptographic Hardware and Embedded Systems, Springer (2009) 156–170