

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

Anonymous Submission

Abstract. In the profiled side-channel analysis, deep learning-based techniques proved to be very successful even when attacking targets protected with countermeasures. Still, there is no guarantee that deep learning attacks will always succeed. What is more, various countermeasures make attacks significantly more difficult, and those countermeasures can be further combined to make the attacks even more challenging. From the other side, to improve the performance of attacks, an intuitive solution would be to reduce the effect of countermeasures.

In this paper, we investigate whether we can consider certain types of hiding countermeasures as noise and then use a deep learning technique called the denoising autoencoder 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.

Keywords: Side-channel analysis, Deep learning, Noise, Countermeasures, Hiding, Denoising autoencoder

1 Introduction

Side-channel analysis (SCA) is a threat exploiting weaknesses in physical implementations of cryptographic algorithms rather than the algorithms themselves [MOP06]. During the execution of an algorithm, leakages like electromagnetic (EM) radiation [QS01] or power dissipation [KJJ99] can happen. Side-channel analysis can be divided into 1) direct attacks like single power analysis (SPA) and differential power analysis (DPA) [KJJ99], and 2) profiled attacks like template attack (TA) [CRR02] and supervised machine learning-based attacks [MPP16, PSK⁺18, CDP17, KPH⁺19]. 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 actively threaten the security of cryptographic devices, there are still severe 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 [BCO04].

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 [CJRR99, BDF⁺17]. 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 [CCD00] or the design of dual-rail logic styles [TV03] are frequently considered options. Exploiting time-randomization is another alternative, e.g., Random Delay Interrupts (RDIs) [CK09] implemented in software and clock jitters in hardware. Still, the countermeasures (especially the hiding ones) are not without weaknesses. Regardless of what hiding approaches are used, we can treat their effects as noise due to their randomness.

44 In other words, the ground truth of the traces always exists. If we can find a way to remove
45 the noise (denoise) from the traces and recover the ground truth of the leakage, then the
46 reconstructed traces could become more vulnerable to the SCA.

47 While considering the countermeasures as noise and then removing that noise sounds like
48 an intuitive approach, this is not an easy problem. The noise (both from the environment
49 and countermeasures) is a part of a signal, and those two components cannot be separated
50 perfectly if we do not know their characterizations. In addition, in realistic settings, we
51 must also consider the portability and the differences among various devices [BCH⁺19].
52 The combination of all these factors makes this problem very complex, and to the best
53 of our knowledge, there are no universal approaches on how to remove the effects of
54 environmental noise and countermeasures.

55 Common approaches to remove/reduce noise are using low-pass filters [WLL⁺18],
56 conducting trace alignments [TGWC18, TGWC18], and various feature engineering meth-
57 ods [ZZY⁺15, PHJB19]. More recently, the SCA community used deep learning techniques
58 that can conduct implicit feature selection and fight countermeasures [CDP17, KPH⁺19,
59 ZBHV19]. While such techniques are useful, they are usually aimed either against a
60 single source of noise or in cases when they can handle more sources of noise, they do
61 not offer interpretability of results. More precisely, in such cases, it is not clear at what
62 point noise removal stops and attack starts (or even if there is such a point). Finally,
63 we emphasize that being able to reduce the noise comprehensively could bring several
64 advantages 1) understanding the attack techniques better, 2) understanding the noise
65 better and consequently, (hopefully) being able to design stronger countermeasures, 3)
66 ability to mount stronger/simpler attacks as there is no noise to consider.

67 In this paper, we propose a new approach to remove several common hiding counter-
68 measures with a denoising autoencoder. Although the denoising autoencoder proved to be
69 successful in removing the noise from an image [Gon16], as far as we are aware, this tech-
70 nique has not been applied to the side-channel domain to reduce the noise/countermeasures
71 effect. We demonstrate the effectiveness of a convolutional denoising autoencoder in dealing
72 with different types of noise and countermeasures separately, i.e., white noise, desynchro-
73 nization, RDIs, and clock jitters. Then, we make the problem more realistic by combining
74 various types of noise and countermeasures with the traces and trying to denoise it with
75 the same machine learning models. The results show that the denoising autoencoder
76 is surprisingly efficient in removing the noise and countermeasures in all investigated
77 situations. We emphasize that denoising autoencoder is not a technique to conduct the
78 profiled attack, but to pre-process the measurements so that any attack strategy can be
79 applied.

80 1.1 Related Work

81 The analysis of the leakage traces in the profiled SCA scenario can be seen as a classification
82 problem where the goal of an attacker is to classify those traces based on the related data
83 (i.e., the encryption key). The most powerful attack from the information-theoretic point
84 of view is the template attack (TA) [CRR02]. Still, this attack can reach its full potential
85 only if the attacker has an unbounded number of traces and the noise follows the Gaussian
86 distribution [LPB⁺15]. More recently, various machine learning techniques emerged as
87 preferred options for cases where 1) the number of traces is either limited or very high,
88 2) the number of features is very high, 3) countermeasures are implemented, and 4) we
89 cannot make assumptions about data distribution. First, the side-channel community
90 showed most interest in techniques like random forest [LMBM13, HPGM16] and support
91 vector machines [HZ12, PHJ⁺17]. More recently, multilayer perceptron [GHO15, PHJ⁺19]
92 and convolutional neural networks [MPP16, CDP17, KPH⁺19] emerged as the most potent
93 approaches. Convolutional neural networks were demonstrated to be capable of coping with
94 the random delay countermeasure due to their spatial invariance property [CDP17, KPH⁺19].

95 At the same time, the fully-connected layers in multilayer perceptron and convolutional
96 neural networks are effective against masking countermeasures as they produce the effect
97 of a higher-order attack (combining features) [BPS⁺18, KPH⁺19]. As far as we are aware,
98 the only application of autoencoders for profiled SCA is made by Maghrebi et al., but
99 there the authors use it for classification and not noise removal, and they report a poor
100 performance when compared to CNNs [MPP16].

101 1.2 Our Contributions

102 In this paper, we consider how to denoise the side-channel traces with convolutional
103 autoencoder (CAE), which to the best of our knowledge, has not been explored before in
104 the side-channel domain. More precisely, we introduce a novel approach to remove the
105 effect of countermeasures and we propose:

- 106 1. A convolutional autoencoder architecture, which requires a limited number of traces
107 to train and can denoise/remove the effect of various hiding countermeasures.
- 108 2. A methodology to recover the ground truth of the traces.

109 To conduct experimental analysis, we consider four separate sources of noise or their com-
110 bination. We investigate the performance of both multilayer perceptron and convolutional
111 neural networks for classification after noise removal. Finally, we experiment with several
112 different noise levels, as well as datasets, having either fixed or random keys, to show the
113 universality of our approach.

114
115 This paper is organized as follows. In Section 2, we provide details about profiled SCAs,
116 machine learning techniques we use, and the dataset we investigate. In Section 3, we provide
117 details about denoising autoencoders and their convolutional version. Section 4 gives
118 experimental results when considering the effects of noise either separately or combined.
119 Finally, in Section 5, we conclude the paper and present possible future research directions.

120 2 Background

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

124 2.1 Notation

125 Let k^* denote the fixed secret cryptographic key (byte), k any possible key hypothesis,
126 and p plaintext. To guess the secret key, the attacker first needs to choose a leakage model
127 $Y(p, k)$ (or Y when there is no ambiguity) depending on the key guess k and some known
128 text p , which relates to the deterministic part of the leakage. The size of the keyspace
129 equals $|K|$. For profiled attacks, the number of acquired traces in the profiling phase equals
130 N , while the number of traces in the testing phase equals T .

131 For the autoencoder, we denote its input as \mathcal{X} . The encoder part of an autoencoder
132 is denoted as ϕ and the decoder part as ψ . Its latent space is denoted as \mathcal{F} . As for the
133 training data, we refer to protected leakages/traces (with noise and countermeasures) as
134 noisy leakages/traces while the unprotected leakages are denoted as clean leakages/traces.

135 2.2 Profiled Side-channel Analysis

136 In the context of implementation attacks, they target physical leakages from the insecure
137 implementation of otherwise theoretically secure cryptographic algorithms. The profiled
138 side-channel attacks represent the most powerful category of SCAs as we assume an

139 attacker with access to an open (keys can be chosen/or are known by the attacker) clone
140 device. Then, the attacker can use that clone device to obtain N measurements from
141 it and construct a characterization model of the device’s behavior. When launching an
142 attack, the attacker then collects a few power traces from the attack device where the
143 secret key is not known. By comparing the attack traces with the characterized model, the
144 secret key can be revealed. Ideally, the secret key can be obtained with a single trace from
145 the attack device. Single trace attack is difficult in practice due to the effect of the noise,
146 countermeasures, and a finite number of traces in the profiling phase (while we assume
147 the attacker is not bounded in his power and he can collect any number of traces, that
148 number represent a small fraction of all possible measurements).

149 To assess the performance of the attacker, one uses a metric denoting the number of
150 measurements required to obtain the secret key. A typical example of such a metric is
151 guessing entropy (GE) [SMY09]. GE represents the average number of key candidates an
152 adversary needs to test to reveal the secret key after conducting a side-channel analysis.
153 In particular, given T amount of samples in the attacking phase, an attack outputs a key
154 guessing vector $g = [g_1, g_2, \dots, g_{|K|}]$ in decreasing order of probability. Then, guessing
155 entropy is the average position of k^* in g over several experiments.

156 2.3 Neural Networks

157 A neural network is an interconnected assembly of simple processing elements, units or
158 nodes, whose functionality is based on the biological process occurring in the brain [Gur14].
159 In general, a neural network consists of three blocks: an input layer, one or more hidden
160 layers, and an output layer, whose processing ability is represented by the strength (weight)
161 of the inter-unit connections, learning from a set of training patterns from the input layer.

162 In the supervised machine learning paradigm, neural networks work in two phases:
163 training and testing. In the training phrases, the goal is to learn a function f , s.t.
164 $f : \mathcal{X} \rightarrow \mathcal{Y}$, given a training set of N pairs (x_i, y_i) . Here, for each example (trace) x , there
165 is a corresponding label y , where $y \in \mathcal{Y}$. Once the function f is obtained, the testing phase
166 starts with the goal to predict the labels for new, previously unseen examples.

167 2.3.1 Multilayer Perceptron

168 The multilayer perceptron (MLP) is a feed-forward neural network that maps sets of inputs
169 onto sets of appropriate outputs. MLP consists of multiple layers (at least three - one
170 input layer, one output layer, and one or more hidden layers) of nodes in a directed graph,
171 where each layer is fully connected to the next one, and training of the network is done
172 with the backpropagation algorithm.

173 2.3.2 Convolutional Neural Networks

174 Convolutional neural networks (CNNs) are a type of neural network originally designed
175 for 2-dimensional convolutions as inspired by the biological processes of animals’ visual
176 cortex [LB+95]. They are commonly used for image classification, but in recent years, they
177 have shown their strengths for time series data [ODZ+16, CFSC17], speech [PC+15], but
178 also security applications [SIJW18].

179 CNNs resemble ordinary neural networks (e.g., multilayer perceptron) from the archi-
180 tecture perspective: they consist of several layers where each layer is made of neurons.
181 CNN usually consists of three types of layers: convolutional layers, pooling layers, and
182 fully-connected layers. Each layer of a network transforms one volume of activation
183 functions to another through a differentiable function. Convolution layer computes the
184 output of neurons that are connected to local regions in the input, each computing a
185 dot product between their weights and a small region they are connected to in the input

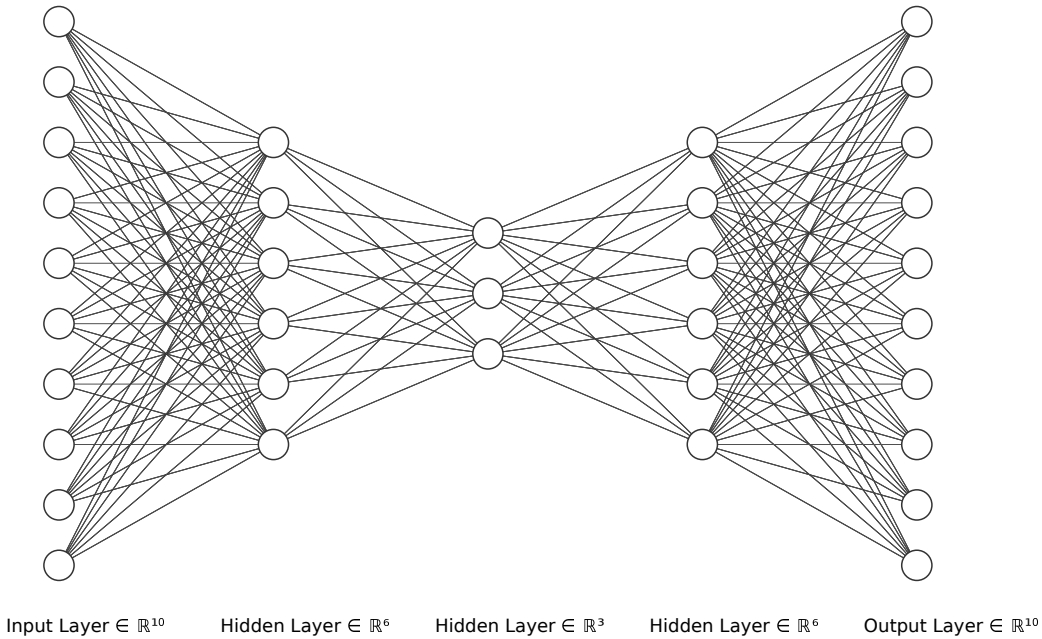


Figure 1: An example of an autoencoder network with three hidden layers (created with NN-SVG [LeN19]).

186 volume. Pooling decrease the number of extracted features by performing a down-sampling
 187 operation along the spatial dimensions. The fully-connected layer computes either the
 188 hidden activations or the class scores. To avoid the overfitting, batch normalization layer,
 189 which normalizes the input layer by adjusting and scaling the activations is commonly
 190 added to the network.

191 2.3.3 Autoencoders

192 Autoencoders were first introduced in the 1980s by Hinton and the PDP group [RHW85] to
 193 address the problem of “backpropagation without a teacher”. Unlike other neural network
 194 architectures that map the relationship between the inputs and the labels, an autoencoder
 195 aims to transform inputs into outputs with the least possible amount of distortion [Bal12].
 196 Benefiting from its unsupervised learning characteristic, autoencoder has been used in
 197 many applications such as data compression [TSC17], anomaly detection [AC15], and
 198 image recovery [Gon16]. The basic structure of an MLP-based autoencoder is given in
 199 Figure 1.

200 An autoencoder consists of two parts: encoder (ϕ) and decoder (ψ). Intuitively, the
 201 encoder squeezes the input with more features to its bottleneck with fewer features, while
 202 the goal of the decoder is to reverse this process. More precisely, the goal of the encoder
 203 is to transfer the input to its latent space \mathcal{F} , i.e., $\phi : \mathcal{X} \rightarrow \mathcal{F}$. The decoder, on the other
 204 hand, reconstructs the input from the latent space, which is equivalent to $\psi : \mathcal{F} \rightarrow \mathcal{X}$.
 205 When training an autoencoder, the goal is to minimize the distortion when transferring
 206 the input to the output (Eq. (1)), i.e., the most representative input features are forced to
 207 be kept in the smallest layer in the network.

$$208 \quad \phi, \psi = \arg \min_{\phi, \psi} \|X - (\psi \circ \phi)X\|^2. \quad (1)$$

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

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

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

224 2.4 The ASCAD Dataset

225 The ASCAD dataset was introduced by Prouff et al. to provide a benchmark to evaluate
 226 machine learning techniques in the context of side-channel attacks [BPS⁺18]. An AT-
 227 Mega8515 device was used to record the emitted EM radiation during the execution of a
 228 software AES implementation protected by known masks. All traces were captured with a
 229 sensor attached to an oscilloscope sampling at 2 GS/s.

230 There are two data sets recorded in different conditions: fixed key encryption and
 231 random key encryption. For the data with fixed key encryption, the dataset provided
 232 separate HDF5 files with different synchronization level: the traces in the ASCAD.h5 file
 233 are time-aligned in a preprocessing step, whereas the traces in ASCAD_desync50.h5 and
 234 ASCAD_desync100.h5 have been shifted with a maximum jitters window of respectively
 235 50 and 100 samples [BPS⁺18]. Each file contains 60 000 EM traces (50 000 training /
 236 cross-validation traces and 10 000 test traces). Each trace consists of 700 points of interest.

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

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

$$244 \quad Y(i) = \text{Sbox}[(p[i] \oplus k[i])]. \quad (2)$$

245 Note that the model given in Eq. (2) does not leak information directly as it is first-order
 246 protected. Consequently, we do not state a model-based SNR. The SNR for the ASCAD
 247 dataset is ≈ 0.8 under the assumption we know the mask (shown in Figure 5a) while it is
 248 almost 0 with the unknown mask.

249 3 Denoising with Convolutional Autoencoder

250 In this section, we discuss the attacker model and the way to obtain clean/noisy traces.
 251 Afterward, we give details about the convolutional autoencoder architecture used in our
 252 experiments.

253 3.1 Denoising Strategy

254 As discussed in Section 2.3.3, noisy-clean trace pairs are required to train a denoising
 255 autoencoder. Inspired by the profiled side-channel analysis method, we devise a denoising
 256 strategy shown in Figure 2.

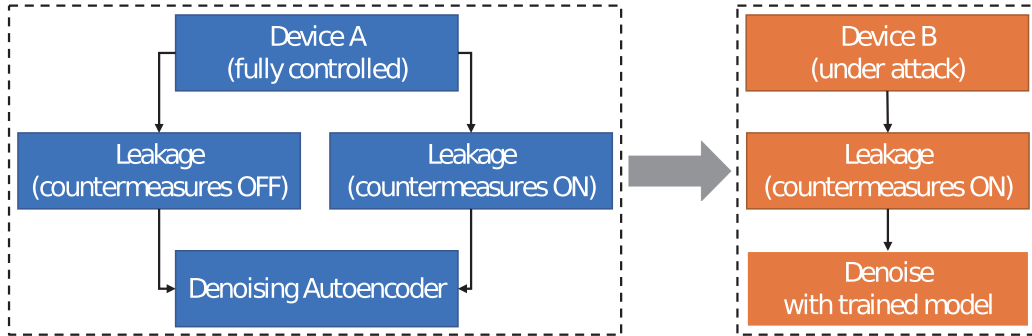


Figure 2: Denoising strategy.

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

265 3.2 Convolutional Autoencoder Architecture

266 An autoencoder can be implemented with different neural network architectures. The most
 267 common examples are the MLP-based autoencoder and convolutional autoencoder (CAE).
 268 Since related works indicate that CNN outperforms MLP in dealing with the side-channel
 269 leakages [MPP16,PHJ⁺19], we use the convolution layer as the basic element for denoising
 270 purpose.

271 To maximize the denoising ability of the proposed architecture, we tune the hyperpa-
 272 rameters by evaluating the CAE performance towards different types of noise. The tuning
 273 range and selected hyperparameters are shown in Table 1. We use the *SeLU* activation
 274 function to avoid vanishing and exploding gradient problems [KUMH17] and *He Uniform*
 275 initialization to improve weight initialization [HZRS15].

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

276 In terms of the autoencoder architecture, we observed that when dealing with trace
 277 desynchronization, an autoencoder that has shallow architecture can denoise the traces
 278 successfully. Still, when introducing other types of noise into the traces while keeping the
 279 same hyperparameters, such autoencoders cannot recover the ground truth of the traces.

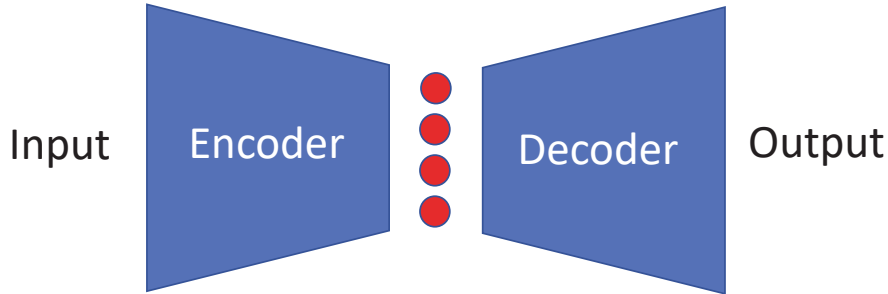


Figure 3: High-level CAE structure

280 Consequently, we decided to increase the depth of the autoencoder to ensure it will be
 281 suitable for different types of noise.

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

291 The details on the CAE architecture used in this paper are given in Table 2. The
 292 convolution block (denoted *Convblock*) usually consists of three parts: convolution layer,
 293 activation layer (function), and Max pooling layer. As we noticed that an autoencoder
 294 implemented in this manner suffers from overfitting and poor performance in denoising
 295 the validation traces, we add the batch normalization layer to each convolution block.

296 A high-level CAE structure is shown in Figure 3. The convolutional encoder converts
 297 the leakages to its latent representation by passing through multiple convolution blocks.
 298 Then, these compressed high-dimensional features are further processed through a fully-
 299 connected layer, which is then reshaped to the convolution size and reconstructed to the
 300 target traces. Note that the size of the latent space (represented by four red circles)
 301 is controlled by the number of neurons in the fully-connected layer. To ensure the CAE
 302 output has the same shape with the training sets, we develop the following equation to
 303 calculate the needed size of the fully-connected layer S_{latent} :

$$304 \quad S_{latent} = \frac{S_{clean}}{\prod_{i=1}^n S_{pooli}} * N_{filter0}. \quad (3)$$

305 S_{clean} is the size of the target clean traces, S_{pooli} represents the i th non-zero pooling stride
 306 of the decoder, and $N_{filter0}$ represents the number of the filters of the first Deconvolution
 307 block. Note, one can vary the size of the latent space for different cases by changing the
 308 size of the pooling layer as well as the number of filters.

309 We emphasize that from an attacker perspective, a CAE can be easily trained by noisy
 310 (protected)-clean (unprotected) traces pairs. Once the training finishes, the autoencoder
 311 can be used to denoise the leakages from real-world devices.

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

Different types of noise must be investigated to evaluate the performance of denoising CAE. At the same time, there are no publicly available datasets concentrating on differences between noise types. To investigate the precise influence of different sources of noise in a fair way, we simulated four types of noise/countermeasures that commonly exist in the devices: Gaussian noise, desynchronization (misalignment), random delay interrupts (RDI), and clock jitters with different noise levels. The simulation approaches are based on previous researches as well as the observation or implementation of the real devices. The pseudocode for the noise simulation is available in the Appendix C.

There is a question of whether the denoising procedure can help regardless of the noise amount. In our experiments, we show the results where the denoising architecture is able to reduce the GE to 0 (or close value) within 5 000 attack traces. Higher noise levels are still affected by CAE, but we require more measurements to reach GE of close to 0. As expected, smaller noise levels are simpler for CAE, and we can reach GE of 0 with even fewer attack traces. We demonstrate this in Appendix A. Finally, we consider what would happen if one applies denoising autoencoder to measurements that do not have noise. In Appendix B, we show that in such a case, a part of the useful signal is removed (as expected), but the attack still works, albeit somewhat worse.

To compare the denoising performance of CAE with the existing techniques commonly used by attackers, we also evaluate the denoising performance of trace averaging. Here, ten 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 an intuition of the properties of the reconstructed traces, we evaluate the SNR of the masked S-box output (Eq. (4)).

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

The SNR is sometimes named F-Test to refer to its original introduction in [Fis22]. For a noisy observation L_t at each time sample t of an event Z , it is defined in Eq. (5). The

340 numerator denotes the signal magnitude and the denominator denotes the noise magnitude
 341 estimation.

$$342 \quad \text{SNR} [t] \triangleq \frac{\text{Var} \mathbb{E}[L_t|Z]}{\mathbb{E} \text{Var}[L_t|Z]}. \quad (5)$$

343 Two models, CNN_best and MLP_best given in the ASCAD paper [BPS⁺18], are
 344 used to attack the traces. The architectures are listed in Table 3 and Table 4. Additionally,
 345 the average pooling layer is used in the Conv block for CNN architecture. In this work,
 346 we use an NVIDIA GTX 1080 Ti graphics processing unit (GPU) with 11 Gigabytes of
 347 GPU memory and 3584 GPU cores. All of the experiments are implemented with the
 348 TensorFlow [AAB⁺15] computing framework and Keras deep learning framework [C⁺15].

Table 3: CNN architecture used for attacking.

Layer	Filter size	Filter number	Pooling stride	Neuron number
Conv block	11	64	2	-
Conv block	11	128	2	-
Conv block	11	256	2	-
Conv block	11	512	2	-
Flatten	-	-	-	-
Fully-connected * 2	-	-	-	4 096

Table 4: MLP architecture used for attacking.

Layer	Neuron number
Fully-connected *4	200

349 Finally, the selected hyperparameters for both deep learning models are shown in
 350 Table 5. Since these two models have different complexities, we can evaluate the denoising
 351 performance in a fair way. We emphasize that we do not aim to find the best attack models
 352 but to show how denoising autoencoders can help improving performance for various
 353 attacks. The quality of the recovered traces is evaluated by guessing entropy (GE). For a
 354 good estimation of GE, the attack traces are randomly shuffled and 100 GEs are computed
 355 to obtain the average value.

Table 5: Hyperparameters for MLP and CNN classifiers.

Hyperparameter	Value
Optimizer	RMSProb
Weight initialization	Uniform distribution
Activation function	ReLU
Learning Rate	1e-5
Batch size	200
Epochs	100 (CNN) / 500 (MLP)
Training sets	35 000
Validation sets	5 000

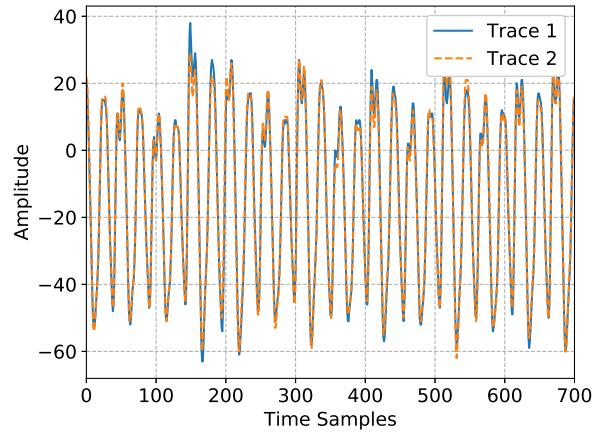
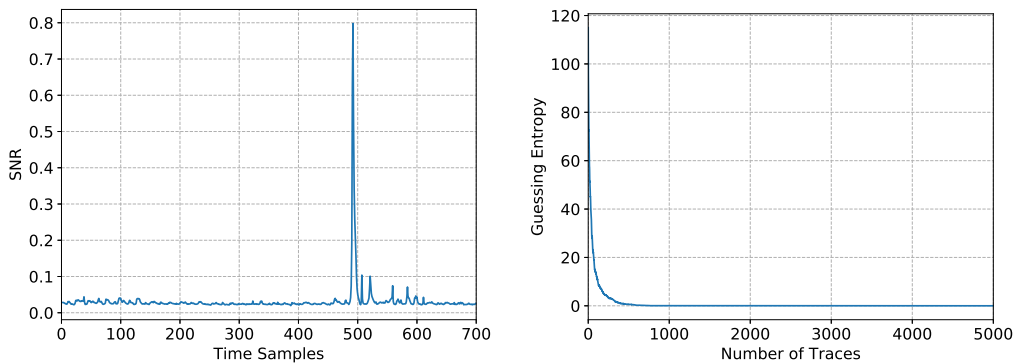


Figure 4: Baseline: example traces.

4.1 Baseline

We first attack the synchronized ASCAD dataset (ASCAD.h5) as the baseline in comparison with different scenarios. The detailed attack approach can be found in the original paper [BPS⁺18]. An example of two traces to be attacked is presented in Figure 4. From the comparison of those two traces, we see that most of the parts are identical, but there are still some variations (especially on the peaks).

Here, 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 5.



(a) SNR: masked S-box output.

(b) Guessing entropy.

Figure 5: Baseline: SNR and guessing entropy.

Without noise and countermeasures in the traces, guessing entropy reaches zero quickly after 405 traces, indicating that the real key has been successfully predicted. The outcomes of this attack are used as the baseline and the following attacks on noisy and denoised traces are compared based on these results.

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,

371 or even the work environment can be the source of Gaussian noise. The noise can also
 372 be artificially introduced by dummy operation or dedicated noise engine. In terms of
 373 trace leakage, the increment of the noise level hides the correlated patterns and reduces
 374 the signal-to-noise (SNR) ratio. Consequently, the noise influences the effectiveness of an
 375 attack, i.e., more traces are needed to obtain the attacked intermediate data.

376 To demonstrate the influence of the Gaussian noise, we add a uniformly distributed
 377 random value between -20 to 20 to each point of the trace. An example of the manipulated
 378 trace and its zoom-in view is shown in Figure 6. Compared with the baseline traces, the
 379 Gaussian noise significantly distorted the shape of the original traces in the amplitude
 380 domain.

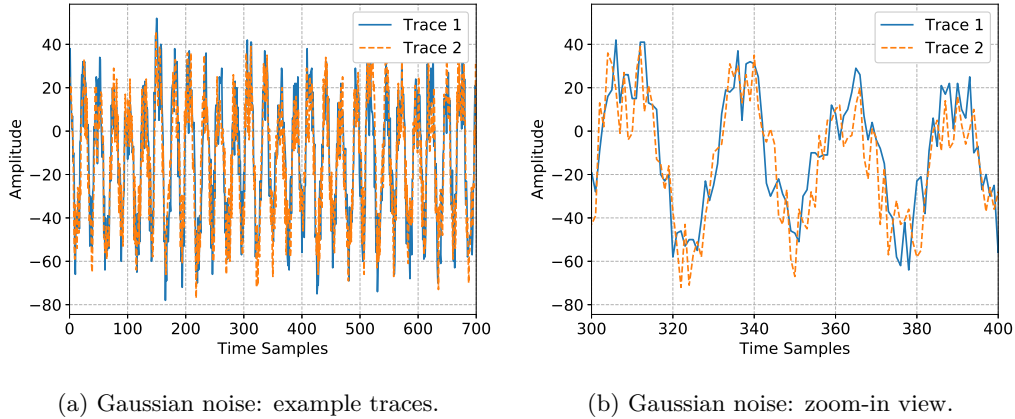


Figure 6: Gaussian noise: example traces and its zoom-in view.

381 Next, we denoise the Gaussian noise with trace averaging as well as CAE proposed in
 382 this paper. The SNR of the noisy and denoised traces are shown in Figure 7a, while GE
 383 values after deep learning attacks are shown in Figure 7b.

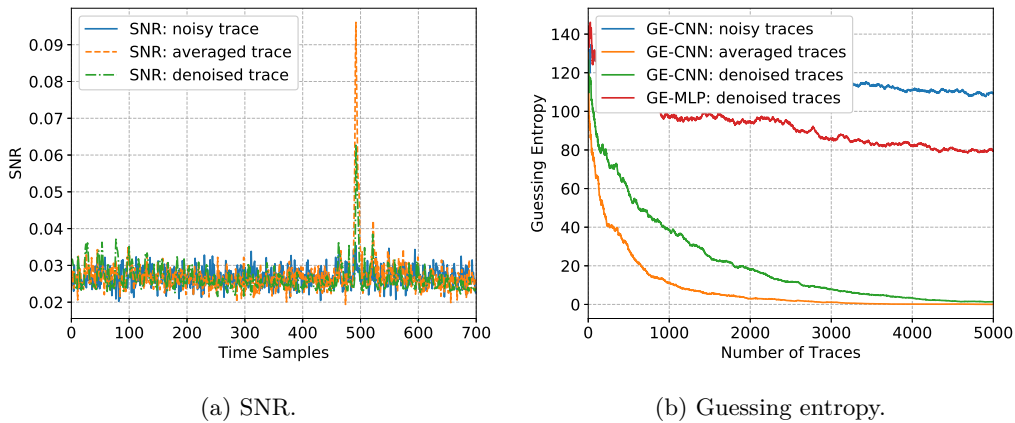


Figure 7: Noisy and denoised traces (Gaussian noise): SNR and guessing entropy.

384 For the baseline attack with “clean” traces, the existence of the Gaussian noise reduces
 385 the leakage of the attacked intermediate data: the SNR drops from 0.8 to 0.03. The SNR
 386 of denoised traces clearly shows that the CAE has filtered out a part of the Gaussian noise,
 387 as the SNR peak of the denoised traces is more than three times higher than the noisy one.

388 From the attack (GE) perspective, GE converges in both cases when the number of

389 trace increases. For the noisy traces, 5 000 traces reduce GE from 121 to only 109 on
 390 average. When considering a multilayer perceptron, we observe that its behavior is better
 391 than CNN when no denoising is done. This indicates that the denoising autoencoder works
 392 regardless of the applied attack technique. What is more, we see that a simpler attack
 393 technique in combination with CAE can outperform more complicated attack techniques.
 394 CNN attack performance after denoising with either averaging or denoising autoencoder
 395 is significantly improved over the noisy version: 5 000 traces are sufficient to reach GE
 396 of 0. It is worth noting that GE of averaged traces is slightly better than GE of CAE,
 397 proving that trace averaging is a good candidate in removing the Gaussian noise. Still,
 398 we can conclude that CAE can remove the Gaussian noise and consequently improve the
 399 attacking efficiency.

400 4.3 Desynchronization

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

412 We consider the desynchronization as a type of noise existing in the traces. Different
 413 from the Gaussian noise, the desynchronization noise adds randomness to the time domain.
 414 To show the effect of the traces desynchronization, we use traces with a maximum of 50
 415 points of desynchronization (ASCAD_desync50 [BPS⁺18]). An example of the traces with
 416 desynchronization is shown in Figure 8.

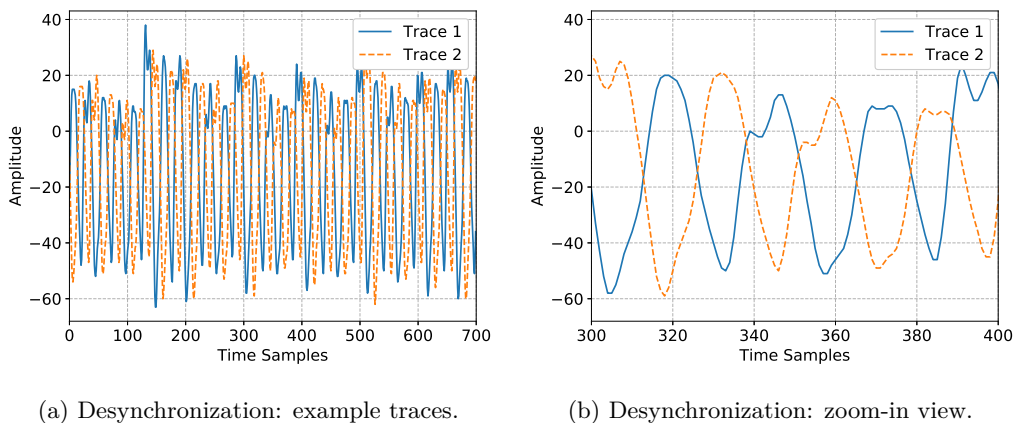


Figure 8: Desynchronization: example traces and its zoom-in view.

417 Next, we attack the misaligned traces as well as the denoised traces from CAE. Figure 9a
 418 shows that only CAE can increase the SNR of the intermediate data (0.16). The averaged
 419 traces, on the other hand, have similar SNR (0.03) with the noisy ones, indicating that
 420 traces averaging is ineffective in dealing with desynchronization.

421 As evident from Figure 9b, GE of the noisy traces converges faster than for the averaged
 422 traces. Consequently, we conclude that the averaging of the desynchronized traces is not

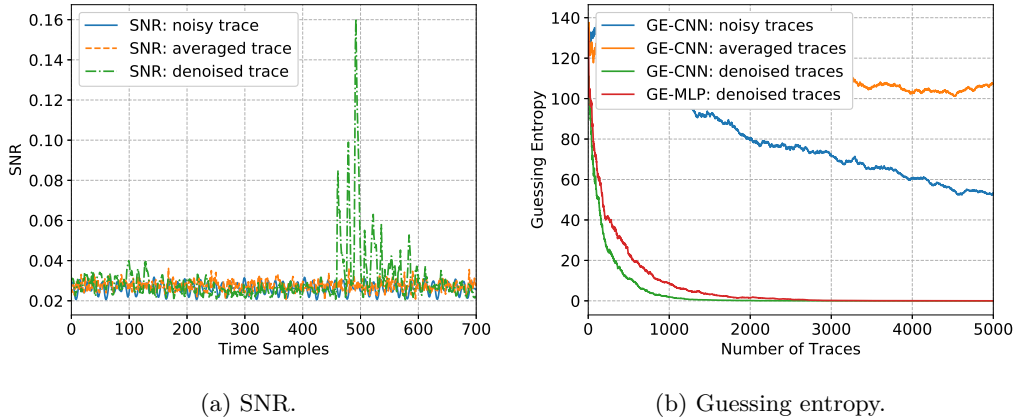


Figure 9: Noisy and denoised traces (desynchronization): SNR and guessing entropy.

423 helping to recover the traces. Indeed, the averaging of the point from different time
 424 locations further degrades the data correlation. On the other hand, CNN proves its ability
 425 to delimit the desynchronization effect, as GE converges to 53 with 5 000 traces when
 426 attacking the noisy traces. Still, considering that the original “clean” traces only needed
 427 405 traces on average to retrieve the key, the desynchronization degraded the performance
 428 of the attack. Additionally, one can expect that performance to become even worse with
 429 an increased desynchronization level.

430 CAE provides a simple alternative approach in synchronizing the traces. By training
 431 a CAE with desynchronized–synchronized traces pair, the model can align the traces
 432 automatically. Consequently, the number of required traces to retrieve the key reduces to
 433 1 189 with CNN and 2 329 with MLP on average.

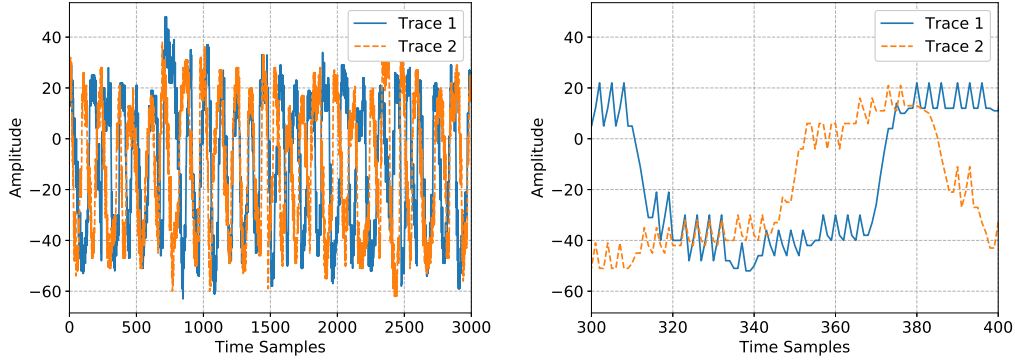
434 4.4 Random Delay Interrupts (RDIs)

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

440 We simulate RDIs based on the Floating Mean method (with parameter $a=5$ and $b=3$)
 441 introduced in [CK09]. The RDIs implemented in such a way can provide more variance
 442 to the traces when compared with the uniform RDI distribution, thus further increasing
 443 the attack difficulty. The probability of the occurrence of RDIs is fixed to 50%. Moreover,
 444 instructions, such as *nop*, are used to generate the random delay according to the real
 445 implementations. In terms of power profile, whenever a random delay occurs, instead of
 446 flattening the power consumption, a specific pattern, such as peak, is shown in the power
 447 trace. We consider this effect by generating a small peak with a certain amplitude when
 448 injecting the random delays to the traces.

449 An example of the traces with RDIs with its zoom-in view is shown in Figure 10. We
 450 observe that more randomness was introduced locally to the traces when compared to the
 451 traces with desynchronization.

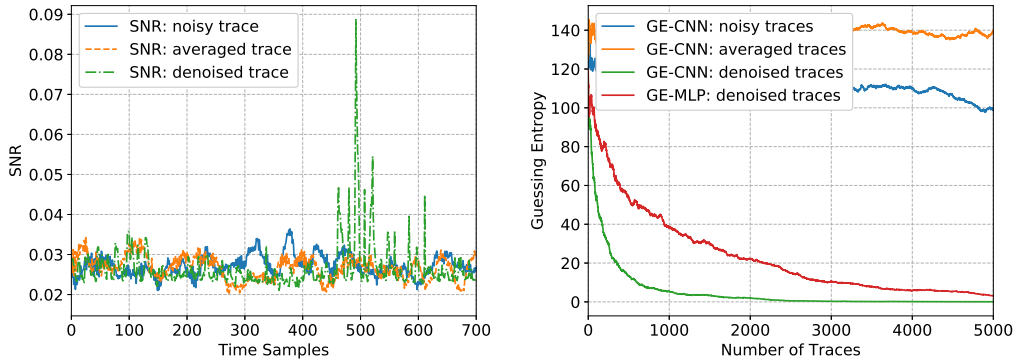
452 As a result, the SNR and rank of the original traces with RDIs and denoised traces
 453 with averaging and CAE are shown in Figure 11. The guessing entropy of the traces with
 454 RDIs converges slowly (129 to 98 with 5 000 traces), indicating that the CNN_best model
 455 we are using is not powerful enough to extract the useful patterns and retrieve the key with



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

Figure 10: Random Delay Interrupts: example traces and its zoom-in view.

456 5 000 traces. We can conclude that RDIs implemented in this way dramatically increase
 457 the attack difficulty.



(a) SNR.

(b) Guessing entropy.

Figure 11: Noisy and denoised traces (RDIs): SNR and guessing entropy.

458 Trace averaging is again not helpful to remove the countermeasure and the GE value
 459 fluctuates around 140. By applying the CAE, the effect of RDIs has been reduced
 460 dramatically: the SNR increases three times to 0.09 and GE converges significantly faster
 461 both when attacking with MLP and CNN. CNN performance is especially good as it needs
 462 only 2 264 traces on average to reach GE of 0 (while MLP requires around double the
 463 amount of traces). We can conclude that CAE can recover the original traces from the
 464 noisy traces with RDIs countermeasure.

465 4.5 Clock Jitters

466 Clock jitters is a classical hardware countermeasure against side-channel attacks, realized
 467 by introducing the instability in the clock [CDP17]. Comparable to the Gaussian noise
 468 that introduces randomness to every point in the amplitude domain, the clock jitters
 469 increase the randomness for each point in the time domain. Indeed, the accumulation of
 470 the deforming effect increases the misalignment of the traces and decreases the correlation
 471 of the intermediate data. As a consequence, the attacked intermediate data become more

472 difficult to retrieve. Here, we simulate the clock jitters by randomly adding or removing
 473 points with a pre-defined range. Similar approaches are used in [CDP17]. More precisely,
 474 we generate a random number r that is uniformly distributed between -4 to 4 to simulate
 475 the clock variation in a magnitude of 8. When scanning each point in the trace, r points
 476 will be added to the trace if r is greater than zero. Otherwise, the following r points in
 477 the trace are deleted. An example of the traces with clock jitters is shown in Figure 12.

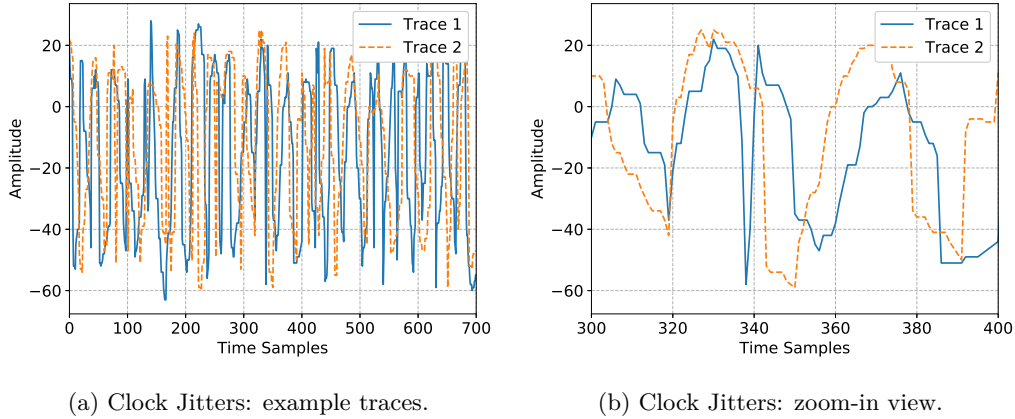


Figure 12: Clock Jitters: example traces and its zoom-in view.

478 We denoise the clock jitters with CAE and a comparison of the attack results for the
 479 noisy and denoised traces with averaging and CAE is shown in Figure 13. As in the
 480 previous settings, traces averaging is not efficient in dealing with time-based noise and
 481 countermeasures. The proposed CAE, on the other hand, successfully reduces the effect of
 482 clock jitters. From the SNR perspective, although the denoised SNR is the lowest (0.044)
 483 for all four different types of noise, it still outperforms its noisy and averaged counterparts.
 484 When compared with noisy traces, the GE of the denoised traces with CNN converges to
 485 13 after applying 5 000 traces¹.

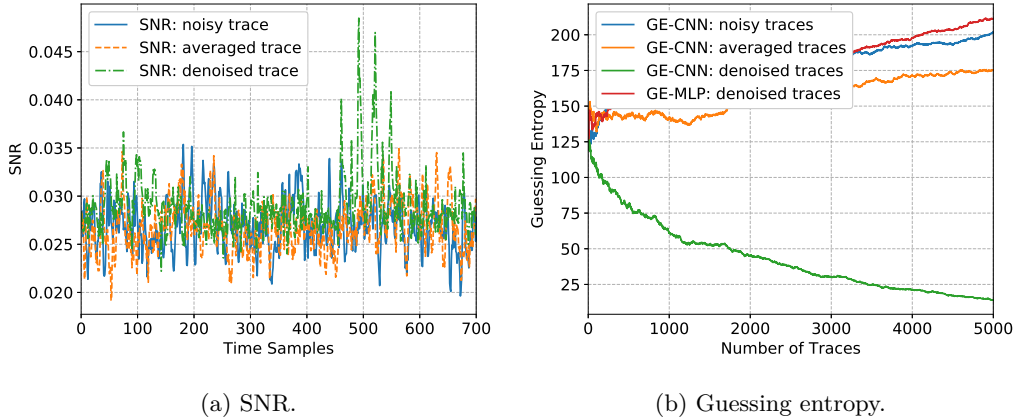


Figure 13: Noisy and denoised traces (clock jitters): SNR and guessing entropy.

486 To conclude, the proposed CAE proves its ability to delimit the effect of the Gaussian

¹Note, we see that for other attacks, GE increases with the increase in the number of traces. We can invert the key guess to obtain the ranking as now the least likely guess would be correct. Still, CNN's performance after DAE is the best one

487 noise, desynchronization, random delay interrupts, and clock jitters. Traces averaging
 488 performs well in removing the Gaussian noise but is ineffective in dealing with the noise
 489 and countermeasures in the time domain. Denoising autoencoder works for both MLP
 490 and CNN attacks, but CNN’s performance is better in comparison.

491 4.6 Combining the Effects of Gaussian Noise and Countermeasures

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

498 4.6.1 AES with the Fixed Key

499 Similar to the procedure of the previous sections, we calculated the SNR and GE of the
 500 noisy and denoised traces and made a comparison between them. The SNR comparison is
 501 presented in Figure 14a. Compared with the noisy traces, the SNR of denoised traces is
 502 increased slightly (0.0345). This observation indicates that the combination of the noise
 503 types degraded the performance of the CAE.

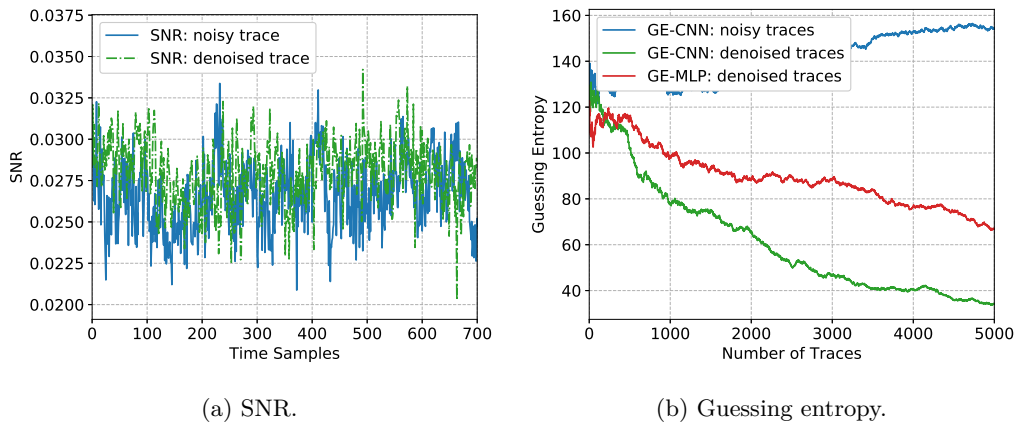


Figure 14: Noisy and denoised traces (all): SNR and guessing entropy.

504 From the GE plot (Figure 14b), the noisy traces do not converge with the increasing
 505 number of traces. The GE of denoised traces, on the other hand, reaches 34 with 5 000
 506 traces when using CNN classifier. Interestingly, the MLP behavior after denoising is
 507 significantly better than CNN’s performance without denoising. Still, GE equals 66 after
 508 applying 5 000 traces. Finally, we can observe that the attack performance converges slower
 509 than the denoised traces with a single type of noise, but CAE still proves its capability in
 510 removing the combined effect of noise and countermeasures.

511 4.6.2 AES with Random Keys

512 Next, we verify the performance of the CAE by trying to denoise the AES traces with
 513 random keys. To retrieve the correct key from the leakage traces, we first train the model
 514 with leakage with a random but known key, then use the trained model to attack the
 515 leakages and try to retrieve the unknown key. When comparing with the fixed-key traces,
 516 the randomness of the key introduces more variance into the traces, thus further increasing

517 the difficulties in denoising the traces. In terms of attack settings, there are 1 400 features
 518 in every trace. The attacked intermediate data was kept the same.

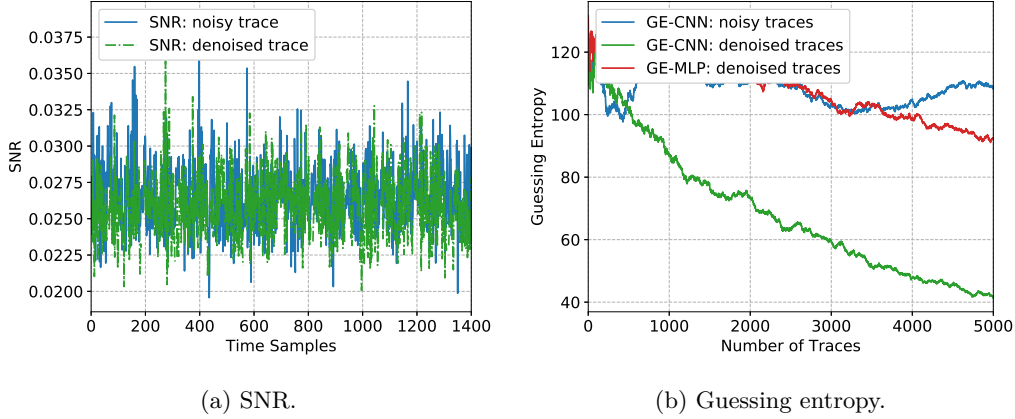


Figure 15: Noisy and denoised traces with variable key (all): SNR and guessing entropy.

519 Based on the results, there is no significant improvement from the SNR perspective.
 520 Still, guessing entropy indicates the improved performance as a result of CAE: for CNN,
 521 GE value converges to 42 with 5 000 traces. The GE of the noisy traces, on the other hand,
 522 fluctuates above 100 regardless of the number of traces. When considering MLP, we see
 523 that denoising helps but reduces GE to only 93 after 5 000 traces. Still, we can conclude
 524 that the proposed CAE can denoise the leakage in both fixed key and variable key scenarios
 525 where the results are especially impressive if using CNN as the attack mechanism.

526 5 Conclusions and Future Work

527 In this paper, we introduce a convolutional autoencoder to remove the noise and counter-
 528 measure from the leakage traces. We consider different types of noise and countermeasures,
 529 such as Gaussian noise, desynchronization, random delay interrupts, and clock jitters.
 530 Additionally, we simulate the scenario where all noise types and countermeasures are com-
 531 bined into the measurements. There, the noisy and denoised traces are compared from two
 532 different perspectives: SNR and guessing entropy. To strengthen our experimental results,
 533 we consider two types of leakage traces (one encrypted with fixed and another with random
 534 keys), two attack strategies (CNN and MLP), and various levels of noise/countermeasure
 535 strengths. The obtained results show that the proposed CAE can still remove/reduce the
 536 noise and find out the underlying ground truth and thus significantly improve the attack
 537 performance.

538 Denoising autoencoder provides an attacker with a powerful tool to pre-process the
 539 traces. Besides that, we expect it could be used to help solve other problems like portabil-
 540 ity [BCH⁺19]. There, the biggest obstacle stems from the variance among different devices.
 541 These variances introduce the variation of the trace, making the attack model generated
 542 for one device difficult to transfer to another one. With the help of an autoencoder, this
 543 problem can be solved in another way: we can consider the traces variation as noise and
 544 use denoising autoencoder to remove it. The basic idea is to select the traces from one
 545 device as a reference and then train an autoencoder to convert the traces from other
 546 devices to the selected device. Since the variation of the traces between devices should
 547 be diminished with the help of the autoencoder, the attack model from one device could
 548 be applied to other devices as it should generalize better. Besides that, in this paper, we

549 considered some common hiding countermeasures. In future work, we aim to investigate
550 whether denoising autoencoder could also work for the masking countermeasures.

551 References

- 552 [AAB⁺15] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen,
553 Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin,
554 Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael
555 Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh
556 Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris
557 Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal
558 Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas,
559 Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and
560 Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous
561 systems, 2015. Software available from tensorflow.org.
- 562 [AC15] Jinwon An and Sungzoon Cho. Variational autoencoder based anomaly detec-
563 tion using reconstruction probability. *Special Lecture on IE*, 2(1), 2015.
- 564 [Bal12] Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In
565 *Proceedings of ICML workshop on unsupervised and transfer learning*, pages
566 37–49, 2012.
- 567 [BCH⁺19] Shivam Bhasin, Anupam Chattopadhyay, Annelie Heuser, Dirmanto Jap,
568 Stjepan Picek, and Ritu Ranjan Shrivastwa. Mind the portability: A warriors
569 guide through realistic profiled side-channel analysis. Cryptology ePrint
570 Archive, Report 2019/661, 2019. <https://eprint.iacr.org/2019/661>.
- 571 [BCO04] Eric Brier, Christophe Clavier, and Francis Olivier. Correlation power analysis
572 with a leakage model. In *International Workshop on Cryptographic Hardware
573 and Embedded Systems*, pages 16–29. Springer, 2004.
- 574 [BDF⁺17] Gilles Barthe, François Dupressoir, Sebastian Faust, Benjamin Grégoire,
575 François-Xavier Standaert, and Pierre-Yves Strub. Parallel implementations
576 of masking schemes and the bounded moment leakage model. In *Annual
577 International Conference on the Theory and Applications of Cryptographic
578 Techniques*, pages 535–566. Springer, 2017.
- 579 [BPS⁺18] Ryad Benadjila, Emmanuel Prouff, Rémi Strullu, Eleonora Cagli, and Cécile
580 Dumas. Study of deep learning techniques for side-channel analysis and
581 introduction to ascad database. *ANSSI, France & CEA, LETI, MINATEC
582 Campus, France. Online verfügbar unter https://eprint.iacr.org/2018/053.
583 pdf, zuletzt geprüft am, 22:2018*, 2018.
- 584 [C⁺15] François Chollet et al. Keras. <https://github.com/fchollet/keras>, 2015.
- 585 [CCD00] Christophe Clavier, Jean-Sébastien Coron, and Nora Dabbous. Differential
586 power analysis in the presence of hardware countermeasures. In *International
587 Workshop on Cryptographic Hardware and Embedded Systems*, pages 252–263.
588 Springer, 2000.
- 589 [CDP17] Eleonora Cagli, Cécile Dumas, and Emmanuel Prouff. Convolutional neural
590 networks with data augmentation against jitter-based countermeasures. In
591 *International Conference on Cryptographic Hardware and Embedded Systems*,
592 pages 45–68. Springer, 2017.

- 593 [CFSC17] Keunwoo Choi, György Fazekas, Mark Sandler, and Kyunghyun Cho. Con-
594 volutional recurrent neural networks for music classification. In *2017 IEEE*
595 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*,
596 pages 2392–2396. IEEE, 2017.
- 597 [CJRR99] Suresh Chari, Charanjit S Jutla, Josyula R Rao, and Pankaj Rohatgi. To-
598 wards sound approaches to counteract power-analysis attacks. In *Annual*
599 *International Cryptology Conference*, pages 398–412. Springer, 1999.
- 600 [CK09] Jean-Sébastien Coron and Ilya Kizhvatov. An Efficient Method for Random
601 Delay Generation in Embedded Software. In *Cryptographic Hardware and*
602 *Embedded Systems - CHES 2009, 11th International Workshop, Lausanne,*
603 *Switzerland, September 6-9, 2009, Proceedings*, pages 156–170, 2009.
- 604 [CRR02] Suresh Chari, Josyula R. Rao, and Pankaj Rohatgi. Template Attacks. In
605 *CHES*, volume 2523 of *LNCS*, pages 13–28. Springer, August 2002. San
606 Francisco Bay (Redwood City), USA.
- 607 [Fis22] Ronald A Fisher. On the mathematical foundations of theoretical statistics.
608 *Philosophical Transactions of the Royal Society of London. Series A, Contain-*
609 *ing Papers of a Mathematical or Physical Character*, 222(594-604):309–368,
610 1922.
- 611 [GHO15] Richard Gilmore, Neil Hanley, and Maire O’Neill. Neural network based attack
612 on a masked implementation of aes. In *2015 IEEE International Symposium*
613 *on Hardware Oriented Security and Trust (HOST)*, pages 106–111. IEEE, 2015.
- 614 [Gon16] Lovedeep Gondara. Medical image denoising using convolutional denoising
615 autoencoders. In *2016 IEEE 16th International Conference on Data Mining*
616 *Workshops (ICDMW)*, pages 241–246. IEEE, 2016.
- 617 [Gur14] Kevin Gurney. *An introduction to neural networks*. CRC press, 2014.
- 618 [HPGM16] Annelie Heuser, Stjepan Picek, Sylvain Guilley, and Nele Mentens. Side-
619 channel analysis of lightweight ciphers: Does lightweight equal easy? In *Radio*
620 *Frequency Identification and IoT Security - 12th International Workshop,*
621 *RFIDSec 2016, Hong Kong, China, November 30 - December 2, 2016, Revised*
622 *Selected Papers*, pages 91–104, 2016.
- 623 [HZ12] Annelie Heuser and Michael Zohner. Intelligent Machine Homicide - Breaking
624 Cryptographic Devices Using Support Vector Machines. In Werner Schindler
625 and Sorin A. Huss, editors, *COSADE*, volume 7275 of *LNCS*, pages 249–264.
626 Springer, 2012.
- 627 [HZRS15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into
628 rectifiers: Surpassing human-level performance on imagenet classification. In
629 *Proceedings of the IEEE international conference on computer vision*, pages
630 1026–1034, 2015.
- 631 [KJJ99] Paul C. Kocher, Joshua Jaffe, and Benjamin Jun. Differential power analysis.
632 In *Proceedings of the 19th Annual International Cryptology Conference on*
633 *Advances in Cryptology*, CRYPTO ’99, pages 388–397, London, UK, UK, 1999.
634 Springer-Verlag.
- 635 [KPH⁺19] Jaehun Kim, Stjepan Picek, Annelie Heuser, Shivam Bhasin, and Alan Hanjalic.
636 Make some noise. unleashing the power of convolutional neural networks for
637 profiled side-channel analysis. *IACR Transactions on Cryptographic Hardware*
638 *and Embedded Systems*, pages 148–179, 2019.

- 639 [KUMH17] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter.
640 Self-normalizing neural networks. In *Advances in neural information processing*
641 *systems*, pages 971–980, 2017.
- 642 [LB⁺95] Yann LeCun, Yoshua Bengio, et al. Convolutional networks for images,
643 speech, and time series. *The handbook of brain theory and neural networks*,
644 3361(10):1995, 1995.
- 645 [LeN19] LeNail. NN-SVG: Publication-Ready Neural Network Architecture Schematics.
646 *Journal of Open Source Software*, 4(33):747, 2019.
- 647 [LMBM13] Liran Lerman, Stephane Fernandes Medeiros, Gianluca Bontempi, and Olivier
648 Markowitch. A Machine Learning Approach Against a Masked AES. In
649 *CARDIS*, Lecture Notes in Computer Science. Springer, November 2013. Berlin,
650 Germany.
- 651 [LPB⁺15] Liran Lerman, Romain Poussier, Gianluca Bontempi, Olivier Markowitch, and
652 François-Xavier Standaert. Template attacks vs. machine learning revisited
653 (and the curse of dimensionality in side-channel analysis). In *International*
654 *Workshop on Constructive Side-Channel Analysis and Secure Design*, pages
655 20–33. Springer, 2015.
- 656 [MOP06] Stefan Mangard, Elisabeth Oswald, and Thomas Popp. *Power Analysis*
657 *Attacks: Revealing the Secrets of Smart Cards*. Springer, December 2006.
658 ISBN 0-387-30857-1, <http://www.dpabook.org/>.
- 659 [MPP16] Housseem Maghrebi, Thibault Portigliatti, and Emmanuel Prouff. Breaking
660 cryptographic implementations using deep learning techniques. In *International*
661 *Conference on Security, Privacy, and Applied Cryptography Engineering*, pages
662 3–26. Springer, 2016.
- 663 [ODZ⁺16] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan,
664 Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray
665 Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint*
666 *arXiv:1609.03499*, 2016.
- 667 [PC⁺15] Dimitri Palaz, Ronan Collobert, et al. Analysis of cnn-based speech recognition
668 system using raw speech as input. Technical report, Idiap, 2015.
- 669 [PHJ⁺17] Stjepan Picek, Annelie Heuser, Alan Jovic, Simone A. Ludwig, Sylvain Guilley,
670 Domagoj Jakobovic, and Nele Mentens. Side-channel analysis and machine
671 learning: A practical perspective. In *2017 International Joint Conference on*
672 *Neural Networks, IJCNN 2017, Anchorage, AK, USA, May 14-19, 2017*, pages
673 4095–4102, 2017.
- 674 [PHJ⁺19] Stjepan Picek, Annelie Heuser, Alan Jovic, Shivam Bhasin, and Francesco
675 Regazzoni. The curse of class imbalance and conflicting metrics with machine
676 learning for side-channel evaluations. *IACR Trans. Cryptogr. Hardw. Embed.*
677 *Syst.*, 2019(1):209–237, 2019.
- 678 [PHJB19] S. Picek, A. Heuser, A. Jovic, and L. Batina. A systematic evaluation of
679 profiling through focused feature selection. *IEEE Transactions on Very Large*
680 *Scale Integration (VLSI) Systems*, 27(12):2802–2815, Dec 2019.
- 681 [PSK⁺18] Stjepan Picek, Ioannis Petros Samiotis, Jaehun Kim, Annelie Heuser, Shivam
682 Bhasin, and Axel Legay. On the performance of convolutional neural networks

- 683 for side-channel analysis. In Anupam Chattopadhyay, Chester Rebeiro, and Yu-
684 val Yarom, editors, *Security, Privacy, and Applied Cryptography Engineering*,
685 pages 157–176, Cham, 2018. Springer International Publishing.
- 686 [QS01] Jean-Jacques Quisquater and David Samyde. Electromagnetic analysis (ema):
687 Measures and counter-measures for smart cards. In Isabelle Attali and Thomas
688 Jensen, editors, *Smart Card Programming and Security*, pages 200–210, Berlin,
689 Heidelberg, 2001. Springer Berlin Heidelberg.
- 690 [RHW85] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning
691 internal representations by error propagation. Technical report, California
692 Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- 693 [SIJW18] Payap Sirinam, Mohsen Imani, Marc Juarez, and Matthew Wright. Deep
694 fingerprinting: Undermining website fingerprinting defenses with deep learn-
695 ing. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and*
696 *Communications Security*, pages 1928–1943. ACM, 2018.
- 697 [SMY09] François-Xavier Standaert, Tal G Malkin, and Moti Yung. A unified framework
698 for the analysis of side-channel key recovery attacks. In *Annual international*
699 *conference on the theory and applications of cryptographic techniques*, pages
700 443–461. Springer, 2009.
- 701 [TGWC18] Hugues Thiebauld, Georges Gagnerot, Antoine Wurcker, and Christophe
702 Clavier. Scatter: A new dimension in side-channel. In *International Workshop*
703 *on Constructive Side-Channel Analysis and Secure Design*, pages 135–152.
704 Springer, 2018.
- 705 [TSCH17] Lucas Theis, Wenzhe Shi, Andrew Cunningham, and Ferenc Huszár.
706 Lossy image compression with compressive autoencoders. *arXiv preprint*
707 *arXiv:1703.00395*, 2017.
- 708 [TV03] Kris Tiri and Ingrid Verbauwhede. Securing encryption algorithms against
709 dpa at the logic level: Next generation smart card technology. In *International*
710 *Workshop on Cryptographic Hardware and Embedded Systems*, pages 125–136.
711 Springer, 2003.
- 712 [WLL⁺18] Lingxiao Wei, Bo Luo, Yu Li, Yannan Liu, and Qiang Xu. I know what you
713 see: Power side-channel attack on convolutional neural network accelerators.
714 In *Proceedings of the 34th Annual Computer Security Applications Conference*,
715 pages 393–406. ACM, 2018.
- 716 [ZBHV19] Gabriel Zaid, Lilian Bossuet, Amaury Habrard, and Alexandre Venelli. Method-
717 ology for efficient cnn architectures in profiling attacks. *IACR Transactions*
718 *on Cryptographic Hardware and Embedded Systems*, 2020(1):1–36, Nov. 2019.
- 719 [ZZY⁺15] Yingxian Zheng, Yongbin Zhou, Zhenmei Yu, Chengyu Hu, and Hailong Zhang.
720 How to compare selections of points of interest for side-channel distinguishers
721 in practice? In Lucas C. K. Hui, S. H. Qing, Elaine Shi, and S. M. Yiu,
722 editors, *Information and Communications Security*, pages 200–214, Cham,
723 2015. Springer International Publishing.

724 A Additional Results with Smaller Levels of Noise

725 To further explore the behavior of the purposed CAE, we explore different levels of
726 noise/countermeasures. For the Gaussian noise, the noise level is reduced from ± 20 to ± 10 .

727 The SNR and GE results are shown in Figure 16. From the results, the GEs of denoised
 728 traces with both MLP and CNN shows better performance than their noisy counterpart.
 729 As expected, since now the noise level is reduced, the GE performance is improved over
 the case with a noise equal to 20.

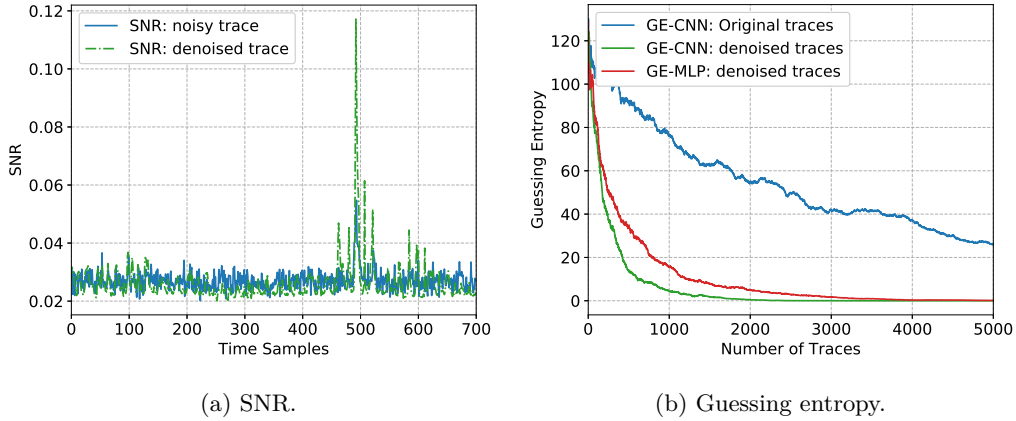


Figure 16: Noisy and denoised traces (Gaussian noise): SNR and guessing entropy.

730

731 For the desynchronization, the noise level is decreased from 50 to 30. The SNR and
 732 GE results are shown in Figure 17. In general, in both cases, fewer traces are needed to
 733 retrieve the key. The GE of the clean traces attacked with CNN is unchanged, indicating
 734 that in both noise levels (high:50/low:30), the maximum denoising capability of CAE is
 reached (also shown in Appendix B).

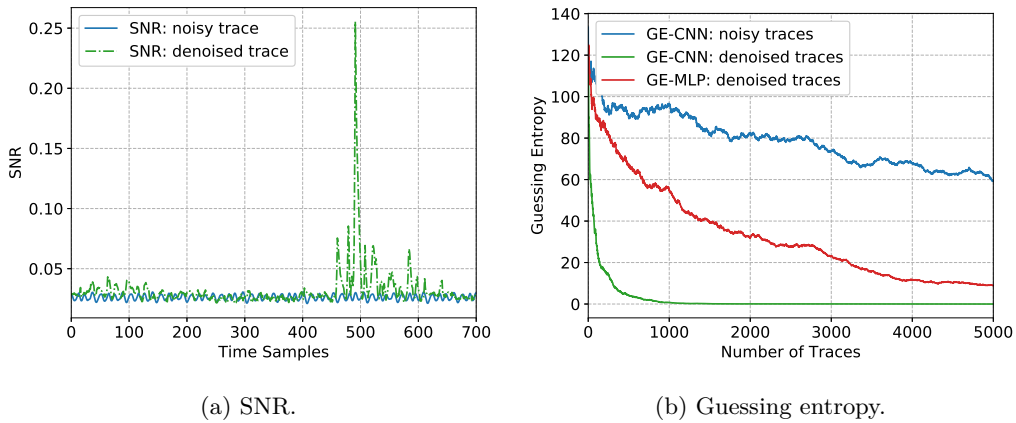


Figure 17: Noisy and denoised traces (desynchronization): SNR and guessing entropy.

735

736 For the random delay interrupts, the noise level is reduced from $a=5/b=3$ to $a=4/b=2$.
 737 The SNR and GE results are shown in Figure 18. As expected, the GE converges faster
 738 for all three attack cases.

739 Finally, for the clock jitters, the noise level was decreased from 4 to 2. The SNR and
 740 GE results are shown in Figure 19. Observe that CNN performs the best and that with
 741 MLP, we need to invert the key guess from the best to the worst one.

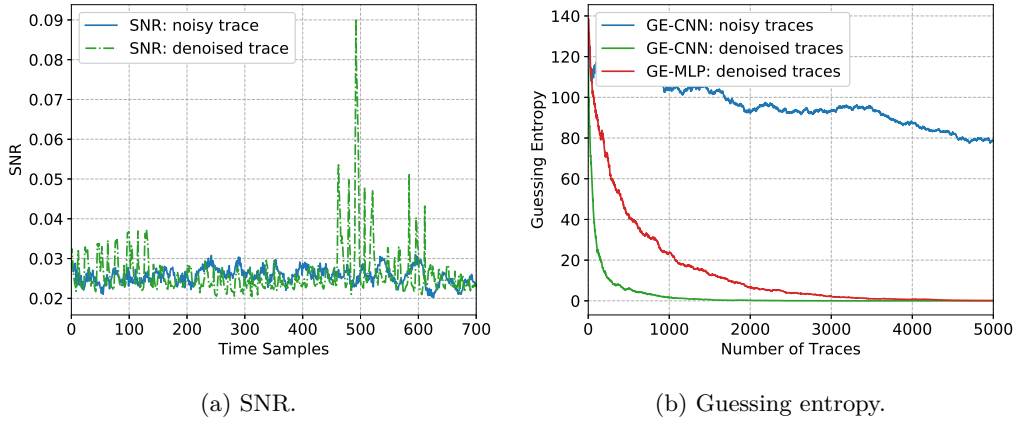


Figure 18: Noisy and denoised traces (RDIs): SNR and guessing entropy.

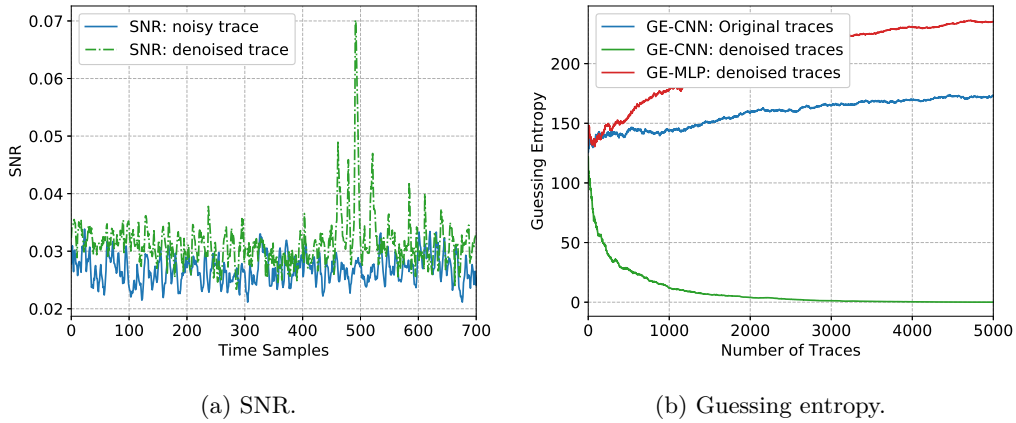


Figure 19: Noisy and denoised traces (clock jitters): SNR and guessing entropy.

B Cleaning the Clean Traces

742

743 One should notice that the traces regenerated by CAE have information loss because of the
 744 bottleneck in the middle of the architecture. In an ideal case, CAE keeps the most critical
 745 information while neglecting the less important features. To evaluate the reconstruction
 746 capability, we now use CAE to denoise the clean traces. In this case, the input and output
 747 of the CAE are identical. While this does not represent a realistic case (as there is no
 748 reason to apply denoising autoencoder to traces that do not have noise), we conduct this
 749 experiment to 1) show that CAE removes mostly noise information, and 2) validate that
 750 even if the evaluator applies by mistake CAE, the performance of the attack will not be
 751 significantly reduced.

752

753 Figure 20 depicts results for the scenario when denoising autoencoder would be applied
 754 to already clean traces (ASCAD dataset with the fixed key). As can be expected, SNR
 755 reduces for those “cleaned” traces as CAE removes some of the information from the signal.
 756 For GE, there is no significant difference if we use clean traces or clean traces after CAE.
 757 In fact, only if the number of traces in the attack set is very limited, there will be a slight
 758 difference in the performance. More precisely, when compared with original traces, there
 759 is a 0.15 drop in terms of SNR; the number of traces to obtain the real key increase from
 405 to 891.

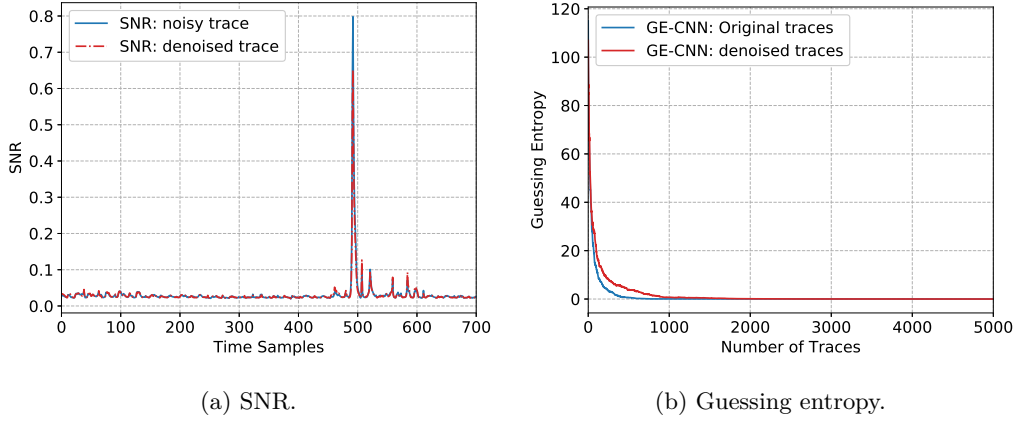


Figure 20: Original clean traces and denoised clean traces: SNR and guessing entropy.

760 C Pseudocode for the Noise Simulation Techniques

761 In Algorithm 1 to 4, we give pseudocode for constructing traces with Gaussian noise,
762 desynchronization, random delay interrupts, and clock jitters, respectively.

Algorithm 1 Gaussian Noise.

```

1: function ADD_GAUSSIAN_NOISE(trace, noise_level)
2:   new_trace  $\leftarrow$  [] ▷ container for new trace
3:   i  $\leftarrow$  0
4:   while i < len(trace) do
5:     level  $\leftarrow$  randomNumber( $-noise\_level$ , noise_level)
6:     new_trace[i]  $\leftarrow$  traces[i] + level ▷ add noise to the trace
7:     i  $\leftarrow$  i + 1
8:   return new_trace

```

Algorithm 2 Desynchronization.

```

1: function ADD_GAUSSIAN_NOISE(trace, desync_level)
2:   new_trace  $\leftarrow$  [] ▷ container for new trace
3:   level  $\leftarrow$  randomNumber(0, desync_level)
4:   i  $\leftarrow$  0
5:   while i + level < len(trace) do
6:     new_trace[i]  $\leftarrow$  traces[i + level] ▷ add desynchronization to the trace
7:     i  $\leftarrow$  i + 1
8:   return new_trace

```

Algorithm 3 Random Delay Interrupts.

```

1: function ADD_RDIS(traces, A, B, rdi_probability, rdi_threshold, rdi_amplitude)
2:   a ← A                                     ▷ maximum number of RDIs
3:   b ← B                                     ▷ a value smaller than A
4:   new_trace ← []                             ▷ container for new trace
5:   i ← 0
6:   while i < len(trace) do
7:     new_trace[i] ← new_trace[i].append(trace[i])
8:     if rdi_probability > rdi_threshold then
9:       m ← randomNumber(0, a - b)
10:      rdi_num ← randomNumber(m, m + b)    ▷ number of RDIs to be added
11:      j ← 0
12:      while j < rdi_num do                 ▷ add RDIs to the trace
13:        new_trace[i] ← new_trace[i].append(trace[i])
14:        new_trace[i] ← new_trace[i].append(trace[i] + rdi_amplitude)
15:        new_trace[i] ← new_trace[i].append(trace[i + 1])
16:        j ← j + 1
17:      i ← i + 1
18:   return new_trace

```

Algorithm 4 Clock Jitters.

```

1: function ADD_CLOCK_JITTERS(trace, clock_jitters_level)
2:   new_trace ← []                             ▷ container for new trace
3:   i ← 0
4:   while i < len(trace) do
5:     new_trace[i] ← new_trace[i].append(trace[i])
6:     level ← randomNumber(0, clock_jitters_level)    ▷ level of clock jitters
7:     if level < 0 then
8:       i ← i + level                             ▷ skip points
9:     else
10:      j ← 0
11:      average_amplitude ← (trace[i] + trace[i + 1])/2
12:      while j < level do
13:        new_trace ← new_trace.append(average_amplitude)    ▷ add points
14:        j ← j + 1
15:      i ← i + 1
16:   return new_trace

```
