

Towards Security Limits in Side-Channel Attacks

(With an Application to Block Ciphers)

F.-X. Standaert*, E. Peeters, C. Archambeau, J.-J. Quisquater
UCL Crypto Group, Place du Levant 3, B-1348 Louvain-la-Neuve, Belgium

Abstract. This paper considers a recently introduced framework for the analysis of physically observable cryptographic devices. It exploits a model of computation that allows quantifying the effect of practically relevant leakage functions with a combination of security and information theoretic metrics. As a result of these metrics, a unified evaluation methodology for side-channel attacks was derived that we illustrate by applying it to an exemplary block cipher implementation. We first consider a Hamming weight leakage function and evaluate the efficiency of two commonly investigated countermeasures, namely noise addition and masking. Then, we show that the proposed methodology allows capturing certain non-trivial intuitions about the respective effectiveness of these countermeasures. Finally, we justify the need of combined metrics for the evaluation, comparison and understanding of side-channel attacks.

1 Introduction

In [14], a unified framework for the analysis of cryptographic primitives against side-channel attacks was introduced as a specialization of Micali and Reyzin’s “physically observable cryptography” paradigm [8]. It exploits a model of computation in which the effect of practically relevant leakage functions is evaluated with a combination of security and information theoretic measurements. A central objective of this framework was to provide a fair evaluation methodology for side-channel attacks. This objective is motivated by the fact that side-channel attacks may take advantage of different statistical tools (*e.g.* difference of means [5], correlation [2], Bayesian classification [1], stochastic models [13]) and are therefore not straightforward to compare. Additionally to the comparisons of side-channel attacks, a more theoretical goal was the understanding of the underlying mechanisms in physically observable cryptography.

Specifically, [14] suggests to combine the success rate of a well specified physical adversary¹ with an information theoretic metric in order to capture the intuition summarized in Figure 1. That is, an information theoretic metric (namely the mutual information) should measure the average amount of information that is available in some physical observations while a security metric measures how efficiently an actual adversary can turn this information into a successful attack, *e.g.* a key recovery that is the goal of most present research on side-channels.

* Postdoctoral researcher funded by the Belgian Fund for Scientific Research (FNRS).

¹ More precisely, [14] suggests either the success rate or the guessing entropy as possible security metrics. This paper only considers the success rate.

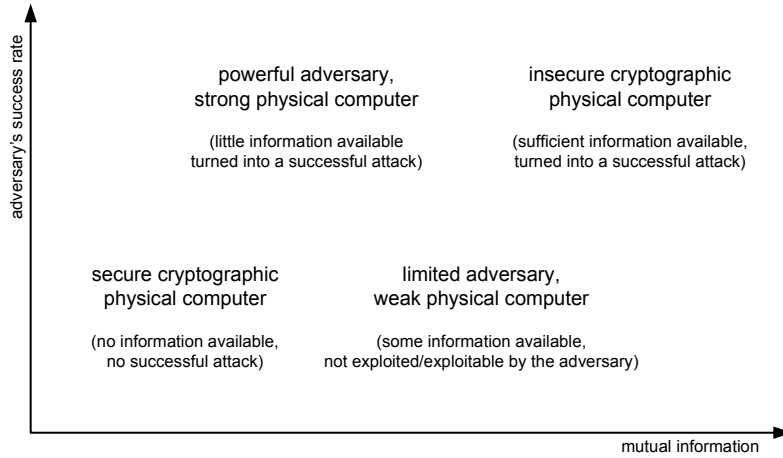


Fig. 1: Summary of side-channel evaluation criteria.

In this paper, we consequently study the relevance of the suggested methodology, by the analysis of a practical case. For this purpose, we investigate an exemplary block cipher and consider a Hamming weight leakage function in different attack scenarios. First, we consider an unprotected implementation and evaluate the information leakages resulting from various number of Hamming weight queries. We discuss how actual block cipher components compare to random oracles with respect to side-channel leakages. Then, we evaluate the security of two commonly admitted countermeasures against side-channel attacks, *i.e.* noise addition and masking. Through these experiments, we show that the proposed evaluation criteria allows capturing certain non-trivial intuitions about the respective effectiveness of these countermeasures. Finally, we provide some experimental validations of our analysis and discuss the advantages of our combination of metrics with respect to other evaluation techniques.

Importantly, in our theoretical framework, side-channel analysis can be viewed as a classification problem. Our results consequently tend to estimate the security limits of side-channel adversaries with two respects. First, because of our information theoretic approach, we aim to evaluate precisely the average amount of information that is available in some physical observations and to determine if this information is sufficient to mount an attack, when an adversary is provided with *unlimited* queries to the device. Second, because we consider (one of) the most efficient classification test(s), namely Bayesian classification, it is expected that the computed success rates also correspond to the best possible adversarial strategy. We mention that the evaluation and comparison metrics to use in the context of side-channel attacks are still under discussion. Our results intend to show that both security and information theoretic metrics are useful, but other similar metrics should still be investigated and compared.

2 Model specifications

In general, the model of computation we consider in this paper is the one initially presented in [8] with the specializations introduced in [14]. In this section, we first describe our target block cipher implementation. Then, we specify the leakage function, adversarial context and decision strategy that we consider in this work. Finally, we provide the definitions of our security and information theoretic metrics for the evaluation of the attacks in the next sections. For a more complete description of the model, we refer to the original paper [14].

2.1 Target implementation

Our target block cipher implementation is represented in Figure 2. For convenience, we only represent the combination of a bitwise key addition and a layer of substitution boxes. We make a distinction between a *single block* and a *multiple block* implementation. This difference refers to the way the key guess is performed by the adversary. In a single block implementation (*e.g.* typically, an 8-bit processor), the adversary is able to guess (and therefore exploit) all the bits in the implementation. In a multiple block implementation (*e.g.* typically, a hardware implementation with data processed in parallel), the adversary is only able to guess the bits at the output of one block of the target design. That is, the other blocks are producing what is frequently referred to as algorithmic noise.

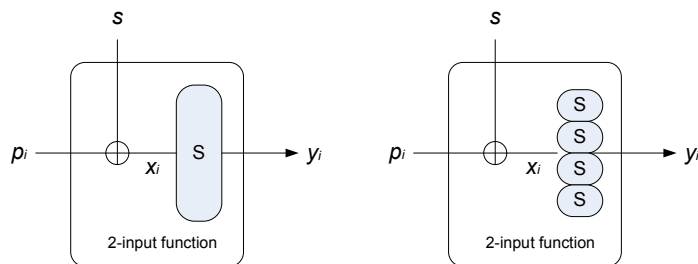


Fig. 2: Single block and multiple block cipher implementation.

2.2 Leakage function

Our results consider the example of a Hamming weight leakage function. Specifically, we assume a side-channel adversary that is provided with the (possibly noisy) Hamming weight leakages of the S-boxes outputs in Figure 2, *i.e.* $W_H(y_i) + n$, where n is a random noise value. As a matter of fact, it involves that our following analysis is theoretical in the sense that we consider simulated leakages. However, since the Hamming weight model has been effectively exploited in a number of works, *e.g.* [2], the obtained conclusions are expected to hold in practice for devices following this type of leakage behavior. Let us finally mention that we only consider univariate leakages with a single leaking point per side-channel query, namely the S-boxes outputs.

2.3 Black box adversarial context and decision strategy

We consider a non-adaptive known plaintext adversary that can perform an arbitrary number of side-channel queries to the target implementation of Figure 2 but cannot choose its queries in function of the previously observed leakages. In addition, we consider a side-channel key recovery adversary with the following (hard) strategy: “*given some physical observations and a resulting classification of key candidates, select the best classified key only*”.

2.4 Security metric: success rate of the key recovery adversary

The success rate of a side-channel key recovery attack can be written as follows. Let S be a discrete variable denoting the target key class in a side-channel attack and s be a realization of this variable. Typically, s corresponds to one or two bytes of the master key. Let \mathbf{L}_q be a random vector denoting the side-channel observations generated with q queries to the target implementation and $\mathbf{l}_q = [l_1, l_2, \dots, l_q]$ be a realization of this random vector, *i.e.* one actual output of the leakage function \mathbf{L} , as defined by Micali and Reyzin in [8]. Following [14], we finally consider a side-channel adversary of which the aim is to guess a key class s with non negligible probability. For this purpose and for each key candidate s^* , it compares the actual observation of a leaking device \mathbf{l}_q with some key dependent model for these leakages $\mathbf{M}(s^*, \cdot)$. Let $\mathbf{T}(\mathbf{l}_q, \mathbf{M}(s^*, \cdot))$ be the statistical test used in the comparison. We assume that the highest value of the statistic corresponds to the most likely key candidate. For each observation \mathbf{l}_q , we first define the set of keys selected by the adversary with a vector:

$$\mathbf{d}_q = \{\hat{s} \mid \mathbf{T}(\mathbf{l}_q, \mathbf{M}(\hat{s}, \cdot)) = \max_{s^*} \mathbf{T}(\mathbf{l}_q, \mathbf{M}(s^*, \cdot))\}.$$

\mathbf{d}_q generally has only one element but several key candidates may have the same test score. Then, we define the result of the attack with the index matrix:

$$\mathbf{I}_{s,s^*}^q = \frac{1}{|\mathbf{d}_q|} \text{ if } s^* \in \mathbf{d}_q, \quad \text{else } 0.$$

Thirdly, we define the success rate of the adversary after q queries:

$$\mathbf{S}_R = \mathbf{E}_s \mathbf{E}_{\mathbf{l}_q|s} \mathbf{I}_{s,s}^q \tag{1}$$

In the following, we will *only* consider a Bayesian classifier, *i.e.* an adversary that selects the keys such that $\Pr[S = s^* | \mathbf{L}_q = \mathbf{l}_q]$ is maximum, since it corresponds to the most efficient way to perform a side-channel key recovery.

It is also interesting to remark that one can use the complete index matrix to build a confusion matrix $\mathbf{C}_{s,s^*}^q = \mathbf{E}_{\mathbf{l}_q|s} \mathbf{I}_{s,s^*}^q$. The previously defined success rate simply corresponds to the averaged diagonal of this matrix. Note finally that the previous definition of success rate corresponds to the success rate against a key class variable (*i.e.* averaged over all possible s) defined in [14], Section 5.1.

2.5 Information theoretic metric: conditional entropy

In addition to the success rate, [14] suggests the use of an information theoretic metric to evaluate the information contained in side-channel observations. Let $\Pr[s|\mathbf{l}_q]$ be the probability of a key candidate s given an observation \mathbf{l}_q with q queries to the target device. We first define an entropy matrix:

$$\mathbf{H}_{s,s^*}^q = - \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s^*|\mathbf{l}_q],$$

from which we derive Shannon's conditional entropy²:

$$\mathbf{H}[S|\mathbf{L}_q] = - \sum_s \Pr[s] \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s|\mathbf{l}_q] = \mathbf{E}_s \mathbf{H}_{s,s}^q \quad (2)$$

We note that this definition is equivalent to the classical one since:

$$\begin{aligned} \mathbf{H}[S|\mathbf{L}_q] &= - \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q] \sum_s \Pr[s|\mathbf{l}_q] \cdot \log_2 \Pr[s|\mathbf{l}_q] \\ &= - \sum_s \Pr[s] \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s|\mathbf{l}_q] \end{aligned}$$

Then, we define an entropy reduction matrix: $\tilde{\mathbf{H}}_{s,s^*}^q = \mathbf{H}[S] - \mathbf{H}_{s,s^*}^q$, where $\mathbf{H}[S]$ is the entropy of the key class variable S before any side-channel attack has been performed: $\mathbf{H}[S] = \mathbf{E}_s - \log_2 \Pr[s]$. It directly yields the mutual information:

$$\mathbf{I}(S; \mathbf{L}_q) = \mathbf{H}[S] - \mathbf{H}[S|\mathbf{L}_q] = \mathbf{E}_s \tilde{\mathbf{H}}_{s,s^*}^q \quad (3)$$

Let us finally mention that in the context of simulated attacks where an analytical model for the leakage probability distribution is known, the previous sums can be turned into integrals, *e.g.* we have for the conditional entropy:

$$\mathbf{H}[S|\mathbf{L}_q] = - \sum_s \Pr[s] \int_{-\infty}^{+\infty} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s|\mathbf{l}_q] d\mathbf{l}_q$$

It is important to observe that the success rate measures the effectiveness of an adversary. In general, it has to be computed for different number of queries in order to evaluate how much observations are required to perform a successful attack. By contrast, the information theoretic metric says nothing about the actual strength of an adversary but characterizes an implementation. It is generally computed once, for an arbitrarily chosen number of queries (typically, $q = 1$).

3 Investigation of single leakages

In this section, we analyze a situation where an adversary is provided with the observation of one single Hamming weight leakage. First, we evaluate single block implementations. Then, we discuss multiple block implementations and key guesses. Finally, we evaluate the effect of noise addition in this context.

² With $\Pr[s|\mathbf{l}_q] = \frac{\Pr[\mathbf{l}_q|s] \cdot \Pr[s]}{\sum_{s^*} \Pr[\mathbf{l}_q|s^*] \cdot \Pr[s^*]}$.

3.1 Single block implementations

Let us assume the following situation: we have an n -bit secret key s and an adversary is provided with the leakage corresponding to a computation $Y_1 = f(s, P_1) = S(P_1 \oplus s)$. As previously, capital letters represent variables while small letters represent particular values of the variables. That is, the adversary obtains observations of the form $I_1 = W_H(y_1)$ and we assume a single block implementation as the one in the left part of Figure 2. Therefore, the adversary can potentially observe the $n + 1$ Hamming weights of the variable Y_1 . Since the Hamming weights of a random value are distributed as binomials, one can easily evaluate the success rate of the adversary as:

$$\mathbf{S}_R = \mathbf{E}_s \mathbf{E}_{I_1} \mathbf{I}_{s,s}^1 = \sum_{h=0}^n \frac{\binom{n}{h}}{2^n} \cdot \frac{1}{\binom{n}{h}} = \frac{n+1}{2^n} \quad (4)$$

This equation means that on average, obtaining the Hamming weight of a secret n -bit value increases the success rate of a key-recovery adversary from $\frac{1}{2^n}$ to $\frac{n+1}{2^n}$. Similar evaluations will be performed for the conditional entropy in Section 3.3.

3.2 Multiple blocks and key guesses

Let us now assume a situation similar to the previous one, but the adversary tries to target a multiple block implementation. Therefore, it is provided with the Hamming weight of an n -bit secret value of which it can only guess b bits, typically corresponding to one block of the implementation. Such a key guess situation can be analyzed by considering the un-exploited bits as a source of algorithmic noise approximated with a Gaussian distribution. This will be done in the next section. The quality of this estimation will then be demonstrated in Section 5, by relaxing the Gaussian estimation.

3.3 Noise addition

Noise is a central issue in side-channel attacks and more generally in any signal processing application. In our specific context, various types of noise are usually considered, including physical noise (*i.e.* produced by the environment), measurement noise (*i.e.* caused by the sampling process and tools), model matching noise (*i.e.* meaning that the leakage model used to attack does possibly not perfectly fit to real observations) or algorithmic noise (*i.e.* produced by the un-targeted values in an implementation). All these disturbances similarly affect the efficiency of a side-channel attack and their consequence is that the information delivered by a single leakage point is reduced. For this reason, a usually accepted method to evaluate the effect of noise is to assume that *there is an additive effect between all the noise sources and their overall effect can be quantified by a Gaussian distribution*. We note that this assumption may not be perfectly verified in practice and that better noise models may allow to improve the efficiency of side-channel attacks. However, this assumption is reasonable in a number of contexts and particularly convenient for a first investigation.

In our experiments, we will consequently assume that the adversary is provided with observations: $\mathbf{l}_{s_g}^1 = W_H(y_1) + n$, where n is a realization of a Gaussian distributed random noise variable N with mean 0 and variance σ^2 . Then, we evaluate the success rate of the adversary and the conditional entropy as:

$$\mathbf{S}_R = \mathbf{E}_s \mathbf{E}_{\mathbf{l}_1} \mathbf{I}_{s,s}^1 = \sum_{h=0}^n \frac{\binom{n}{h}}{2^n} \cdot \int_{-\infty}^{+\infty} \Pr[\mathbf{l}_1|h] \cdot \mathbf{I}_{s,s}^1 dl, \quad (5)$$

$$\mathbf{H}[S|\mathbf{L}_1] = \mathbf{E}_s \mathbf{H}_{s,s}^1 = \sum_{h=0}^n \frac{\binom{n}{h}}{2^n} \cdot \int_{-\infty}^{+\infty} \Pr[\mathbf{l}_1|h] \cdot -\log_2(\Pr[s|\mathbf{l}_1]) dl, \quad (6)$$

where $\Pr[\mathbf{L}_1 = \mathbf{l}_1 | W_H(Y_1) = h] = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\mathbf{l}_1 - h)^2}{2\sigma^2}\right)$ and the a posteriori probability $\Pr[s|\mathbf{l}_1]$ can be computed thanks to Bayes's formula: $\Pr[s|\mathbf{l}_1] = \frac{\Pr[\mathbf{l}_1|s] \cdot \Pr[s]}{\Pr[\mathbf{l}_1]}$, with $\Pr[\mathbf{l}_1] = \sum_{s^*} \Pr[\mathbf{l}_1|s^*] \cdot \Pr[s^*]$. As an illustration, the success rate and the mutual information are represented in Figure 3 for an 8-bit value, in function of the observation signal-to-noise ratio ($\text{SNR} = 10 \cdot \log_{10}\left(\frac{\varepsilon^2}{\sigma^2}\right)$, where $\varepsilon = \sqrt{n}/4$ denotes the standard deviation of the Hamming weight signal and σ is the previously introduced Gaussian noise standard deviation).

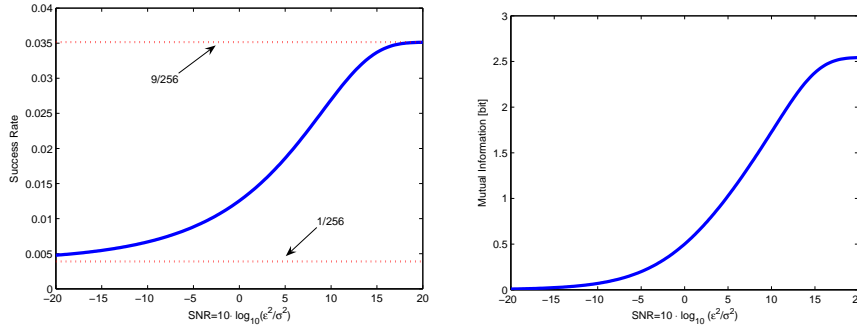


Fig. 3: Success rate and mutual information in function of the SNR.

Note that the success rate starts at $9/256$, *i.e.* the noise-free value computed with Equation (4) and tends to $1/256$ which basically means that very little information can be retrieved from the leakage. Also, since for each key class s , the same 9 Hamming weights can be observed with the same frequency, we are typically in the context of a weak template attack as described in [14]. That is, each line of the entropy matrix is identical up to a permutation of its elements.

4 Investigation of multiple leakages

In the previous section, we analyzed a single-query adversary. However, looking at Figure 3, it is clear that such a context involves limited success rates, even in case of high SNRs. As a matter of fact, actual adversaries would not only perform one single query to the target device but multiple ones, in order to increase their success rates. Therefore, this section considers the problem of multiple leakages.

For this purpose, let us consider the following situation: we have an n -bit secret key class s and an adversary is provided with the leakages corresponding to two computations $Y_1 = f(s, P_1)$ and $Y_2 = f(s, P_2)$. That is, it obtains $W_H(Y_1)$ and $W_H(Y_2)$ and we would like to evaluate the average predictability of s . The

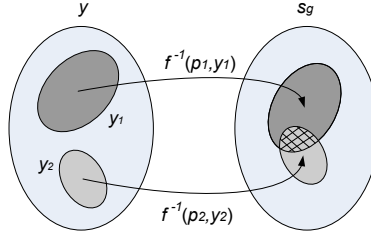


Fig. 4: Multiple point leakages.

consequence of such an experiment (illustrated in Figure 4) is that the key will be contained in the intersection of two sets of candidates obtained by inverting the 2-input functions $y_1 = f(s, p_1)$ and $y_2 = f(s, p_2)$. The aim of our analysis is therefore to determine how the keys within this intersection are distributed. Importantly, and contrary to the single query context, this analysis requires to characterize the cryptographic functions used in the target implementation, since they will determine how the intersection between the sets of candidates behaves. Therefore, we will consider two possible models for these functions.

4.1 Assuming random S-boxes

A first (approximated) solution is to consider the functions $f^{-1}(P_i, Y_i)$ to behave randomly. As a consequence, each observed Hamming weight leakage $h_i = W_H(y_i)$ will give rise to a uniform list of candidates for the key s of size $n_i = \binom{n}{h_i}$, without any particular dependencies between these sets but the key. Let us denote the size of the set containing s after the observation of q leakages respectively giving rise to these uniform lists of n_i candidates by a random variable $T_q(n_1, n_2, \dots, n_q)$. From the probability density function of T_q (given in appendix A), it is straightforward to extend the single leakage analysis of Section 3.1 to multiple leakages. The success rate can be expressed as:

$$\mathbf{S}_R = \sum_{h_1=0}^n \sum_{h_2=0}^n \dots \sum_{h_q=0}^n \frac{\binom{n}{h_1}}{2^n} \cdot \frac{\binom{n}{h_2}}{2^n} \dots \frac{\binom{n}{h_q}}{2^n} \cdot \sum_i \Pr[T_q = i] \cdot \frac{1}{i} \quad (7)$$

4.2 Using real block cipher components

In order to validate the previous theoretical predictions of the success rate, we performed the experiments illustrated in Figure 5. In the first (upper) experiment, we generated a number of plaintexts, observed the outputs of the function $f = S(P_i \oplus s)$ through its Hamming weights $W_H(Y_i)$, derived lists of n_i candidates for Y_i corresponding to these Hamming weights and went through the

inverted function $f^{-1}(P_i, Y_i)$ to obtain lists of key candidates. In the second (lower) experiment, a similar procedure is applied but the n_i key candidates were selected from random lists (including the correct key). As a matter of fact, the first experiment corresponds to a side-channel attack against a real block cipher (we used the AES Rijndael S-box) while the second experiment emulates the previous random S-box estimation.

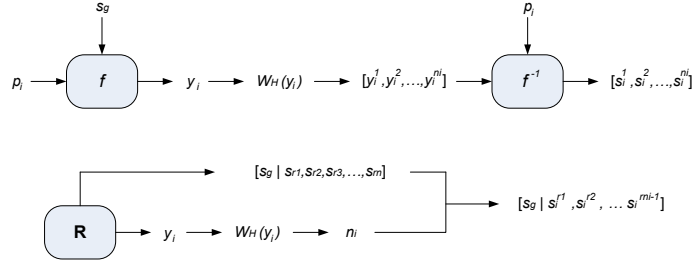


Fig. 5: Multiple leakages experiments: real S-boxes and random S-boxes simulation.

We generated a large number (namely 100 000) of observations and, for these generated observations, derived the experimental success rate in the two previous contexts. Additionally, we compared these experiments with the theoretical predictions of the previous section. The results of our analysis are pictured in Figure 6, where we can observe that the real S-box gives rise to lower success rates (*i.e.* to less information) than a random function. The reason of this phenomenon is that actual S-boxes give rise to (slightly) correlated lists of key candidates and therefore to less independence between consecutive observations, as already suggested in [2, 11]. These experiments suggest that even if not perfectly correct, the assumption that block cipher components are reasonably approximated by random functions with respect to side-channel attacks is acceptable. We note that this assumption is better verified for large bit sizes since large S-boxes better approximate the behavior of a random function than small ones.

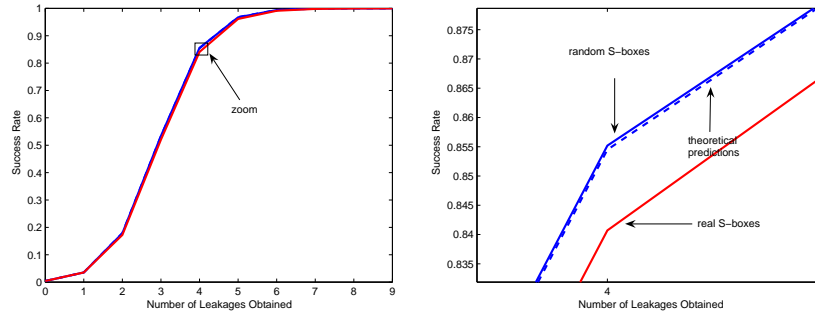


Fig. 6: Multiple leakages experimental results.

5 Investigation of masked implementations

The previous sections illustrated the evaluation of simple side-channel attacks based on a Hamming weight leakage function thanks to the success rate and mutual information. However, due to the simplicity of the investigated contexts, these notions appeared to be closely correlated. Therefore it was not clear how one could need both criteria for our evaluation purposes. In this section, we consequently study a more complex case, namely masked implementations and higher-order side-channel attacks. This example is of particular interest since it allows us to emphasize the importance of a combination of security and information theoretic metrics for the physical security evaluation process of an implementation. As a result of our analysis, we provide (non-trivial) observations about the respective effectiveness of masking and algorithmic noise addition that can be easily turned into design criteria for actual countermeasures.

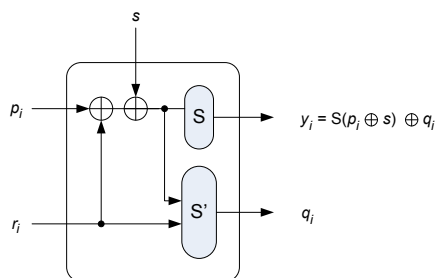


Fig. 7: 1st order boolean masking.

The masking technique (*e.g.* [4]) is one of the most popular ways to prevent block cipher implementations from Differential Power Analysis. However, recent results suggested that it is not as secure as initially thought. Originally proposed by Messerges [7], second and higher-order power analysis attacks can be successfully implemented against various kinds of designs and may not require more hypotheses than a standard DPA [9]. In [12], an analysis of higher-order masking schemes is performed with respect to the correlation coefficient. In the following, we intend to extend this analysis to the (more powerful but less flexible) case of a Bayesian adversary, as introduced in [10].

For the purposes of our analysis, we will use the masked implementation illustrated in Figure 7 in which the plaintext p_i is initially XORed with a random mask r_i . We use two S-boxes S and S' such that: $S(p_i \oplus r_i \oplus s) = S(p_i \oplus s) \oplus q_i$, with $q_i = S'(p_i \oplus r_i \oplus s, r_i)$. According to the notations introduced in [10], it is particularly convenient to introduce the secret state of the implementation as $\Sigma_i = S(p_i \oplus s)$ and assume an adversary that obtains (possibly noisy) observations: $\mathbf{l}_q = W_H[\Sigma_i \oplus q_i] + W_H[q_i] + n$, with the same noise as in Section 3.3. Similarly to a first-order side-channel attack, the objective of an adversary is then to determine the secret state Σ_i (it directly yields the secret key class s). Because of the masking, Σ_i is not directly observable through side-channel measurements but its associated PDFs do, since these PDFs only depend on

the Hamming weight of the secret state $W_H(\Sigma_i)$. As an illustration, we provide the different discrete PDFs (computed over the random mask values) for a 4-bit masked design in Figure 8, in function of the secret state Σ_i . We also depict the shapes of the discrete PDFs corresponding to an unmasked secret state affected by four bits of algorithmic noise (*i.e.* we add 4 random bits to the 4-bit target and the PDF is computed over these random bits). Similar distributions can be obtained for any bit size. In general, knowing the probability distributions of the secret state, the success rate and conditional entropy can be straightforwardly derived. For example, after one query it yields:

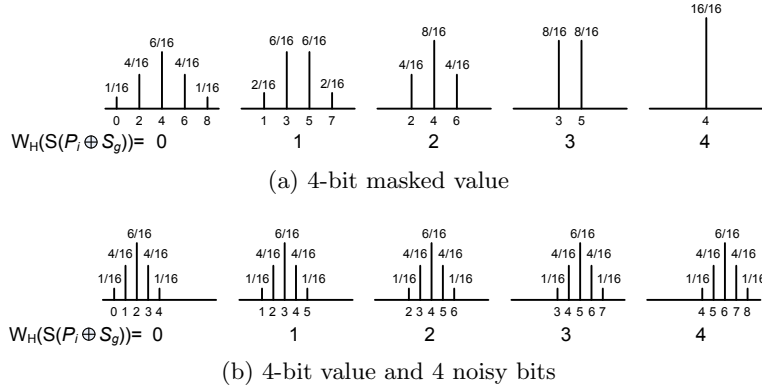


Fig. 8: Exemplary discrete leakage PDFs.

$$\mathbf{S}_R = \mathbf{E}_{\Sigma_1} \mathbf{E}_{\mathbf{I}_1} \mathbf{I}_{\Sigma_1, \Sigma_1}^1 = \sum_{h=0}^n \frac{\binom{n}{h}}{2^n} \cdot \int_{-\infty}^{+\infty} \Pr[\mathbf{I}_1|h] \cdot \mathbf{I}_{\Sigma_1, \Sigma_1}^1 dl, \quad (8)$$

$$\mathbf{H}[S|\mathbf{L}_1] = \mathbf{E}_{\Sigma_1} \mathbf{H}_{\Sigma_1, \Sigma_1}^1 = \sum_{h=0}^n \frac{\binom{n}{h}}{2^n} \cdot \int_{-\infty}^{+\infty} \Pr[\mathbf{I}_1|h] \cdot -\log_2(\mathbf{P}[\Sigma_1|\mathbf{I}_1]) dl, \quad (9)$$

where $\Pr[\mathbf{L}_1 = \mathbf{l}_1 | W_H(\Sigma_1) = h]$ can be computed as in Section 3.3, assuming that the \mathbf{I}_1 's are distributed as a mixture of Gaussians. In the following, we illustrate these metrics in different contexts. First, we consider 1st and 2nd order masking schemes for 8-bit S-boxes. Then, we consider unmasked implementations where 8 (*resp.* 16) random bits of algorithmic noise are added to the secret signal S , corresponding to the 1st (*resp.* 2nd) order mask bits.

The first (and somewhat surprising) conclusion of our experiments appears in Figure 9. Namely, looking at the mutual information for high SNRs, the use of a n -bit mask is less resistant (*i.e.* leads to lower leakages) than the addition of n random bits to the implementation. Fortunately, beyond a certain amount of Gaussian noise the masking appears to be a more efficient protection. The reason of this behavior appears clearly when observing the evolution of the PDFs associated to each secret state in function of the SNR, pictured in Appendix B, Figures 13 and 14. Clearly, the PDFs of the masked implementation are very

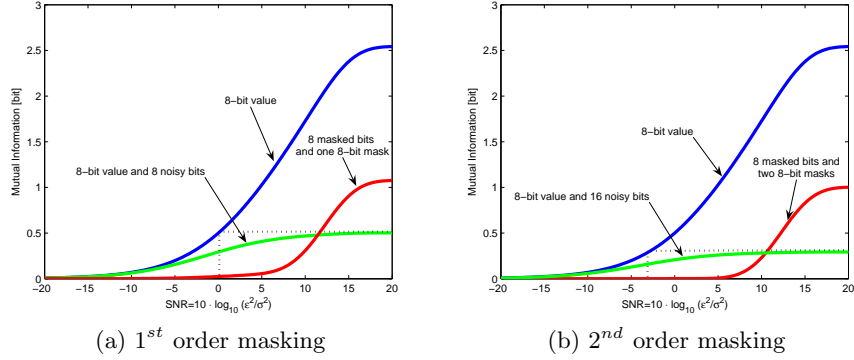


Fig. 9: Mutual information of 1^{st} , 2^{nd} order masking and equivalent algorithmic noise.

different with small noise values (*e.g.* in Figure 13.a, the probability that an observation belong to both PDFs is very small) but becomes almost identical when the noise increases, since they are all identically centered (*e.g.* in Figure 13.b). Conversely, the means of each PDF in the unmasked implementations stay different whatever the noise level (*e.g.* in Figure 14.b). Therefore the Bayesian classification is easier than in the masked case when noise increases. These observations confirm the usually accepted fact that efficient protections against side-channel attacks require to combine different countermeasures. A practically important consequence of our results is the possibility to derive the exact design criteria (*e.g.* the required amount of noise) to obtain an efficient masking.

It is also interesting to observe that Figure 9 confirms that algorithmic noise is nicely modeled by Gaussians. Indeed, *e.g.* for the 1^{st} order case, the mutual information of an 8-bit value with 8 noisy bits for high SNRs exactly corresponds to the one of an unprotected 8-bit value with $SRN=0$.

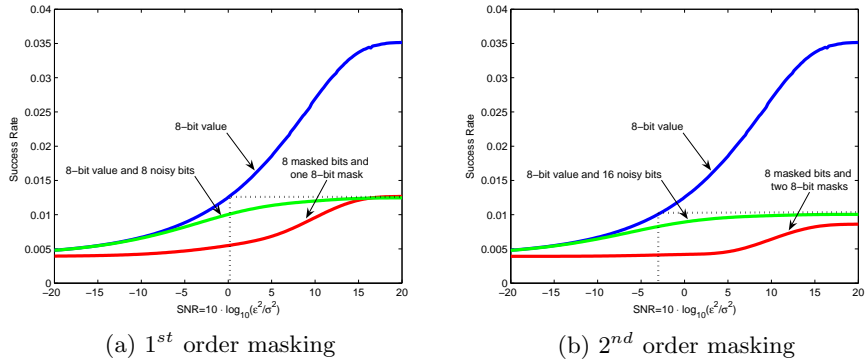


Fig. 10: Success rate of 1^{st} , 2^{nd} order masking and equivalent algorithmic noise.

The second interesting conclusion is that the success rate after one query (pictured in Figure 10) does *not* follow an identical trend. Namely, the masked implementations and their equivalent noisy counterparts do *not* cross over at the same SRN. This situation typically corresponds to the intuitive category

of limited adversary against a weak implementation in Figure 1. That is, some information is available but the number of queries is too low to turn it into a successful attack. If our information theoretic measurement is meaningful, higher number of queries should therefore confirm the intuition of Figure 9.

Success rates with higher number of queries for a 2^{nd} order masking scheme (and noisy equivalent) were simulated in Figures 11, 12. In Figure 11, a very high SNR=20 is considered. As a consequence, we observe that the masks bring much less protection than their equivalent in random bits, although the initial value (for one single query) suggests the opposite. Figure 12 performs similar experiments for two SNRs that are just next to the crossing point. It illustrates the same intuition that the efficiency of the key recovery when increasing the number of queries is actually dependent on the information content in the observations.

Importantly, these experiments illustrate a typical context where the combination of security and information theoretic metrics is meaningful. While the success rate is the only possible metric for the comparison of different side-channel attacks (since it could be evaluated for different statistical tools), the information theoretic metric allows to infer the behavior of an attack when increasing the number of queries. As an illustration, the correlation-based analysis performed in [12] only relates to one particular (sub-optimal) statistical tool and was not able to lead to the observations illustrated in Figure 9.

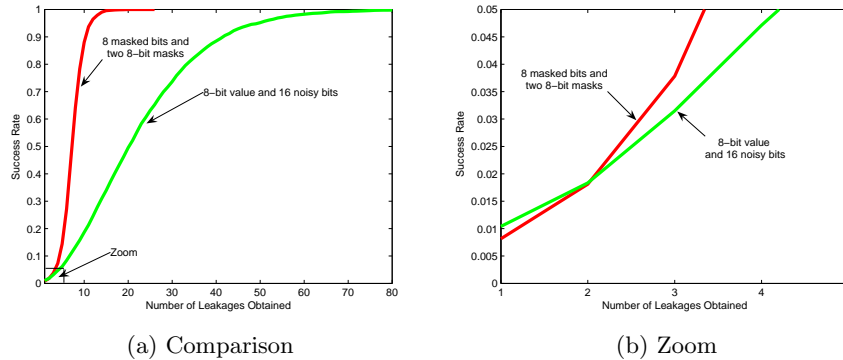


Fig. 11: Success rate of an 8-bit 2^{nd} order masking scheme with noisy counterpart.

6 Concluding remarks

This paper discusses the relevance of a recently introduced theoretical framework for the analysis of cryptographic implementations against side-channel attacks. By the investigation of a number of implementation contexts, we illustrate the interest of a combination of security and information theoretic metrics in the evaluation, comparison and understanding of side-channel attacks. Specifically, our results show a practically meaningful example in which computing the mutual information of the leakages provides theoretical insights about the asymptotic security of an implementation that can possibly be turned into practical

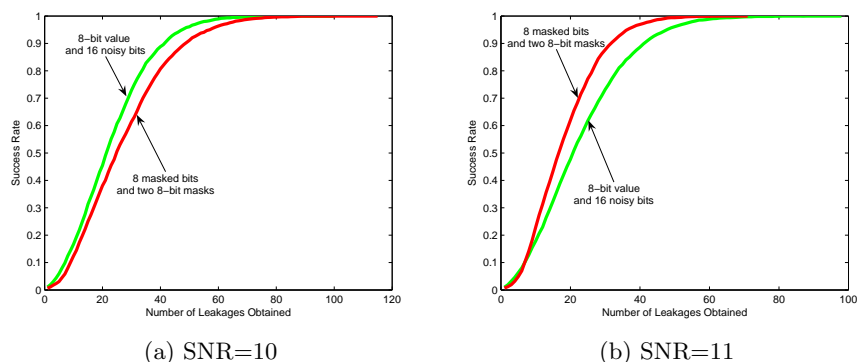


Fig. 12: Success rate of an 8-bit 2^{nd} order masking scheme with noisy counterpart.

design criteria. As a scope for further research, we suggest the analysis of more complex (statistically sampled, multivariate, ...) leakage functions possibly giving rise to strong template attacks (as defined in [14]), *i.e.* attacks in which each key class gives rise to a different security level.

References

1. S. Chari, J.R. Rao, P. Rohatgi, *Template Attacks*, CHES 2002, Lecture Notes in Computer Science, vol. 1965, pp. 13–28.
2. E. Brier, C. Clavier, F. Olivier, *Correlation Power Analysis with a Leakage Model*, CHES 2004, Lecture Notes in Computer Science, vol. 3156, pp. 16–29.
3. J.-S. Coron, P. Kocher, D. Naccache, *Statistics and Secret Leakage*, Financial Crypto 2000, Lecture Notes in Computer Science, vol. 1972, pp. 157–173.
4. L. Goubin, J. Patarin, *DES and Differential Power Analysis*, CHES 1999, Lecture Notes in Computer Science, vol. 1717, pp. 158–172.
5. P. Kocher, J. Jaffe, B. Jun, *Differential Power Analysis*, CRYPTO 1999, Lecture Notes in Computer Science, vol. 1666, pp. 15–19.
6. S. Mangard, *Hardware Countermeasures against DPA - a Statistical Analysis of their Effectiveness*, CT-RSA 2004, LNCS, vol. 2964, pp. 222–235.
7. T.S. Messerges, *Using Second-Order Power Analysis to Attack DPA Resistant Software.*, CHES 2000, LNCS, vol. 2523, pp. 238–251.
8. S. Micali, L. Reyzin, *Physically Observable Cryptography (extended abstract).*, TCC 2004, Lecture Notes in Computer Science, vol. 2951, pp. 278–296.
9. E. Oswald, S. Mangard, C. Herbst, S. Tillich, *Practical Second-Order DPA Attacks for Masked Smart Card Implementations of Block Ciphers*, CT-RSA 2006, Lecture Notes in Computer Science, vol. 3860, pp. 192–207.
10. E. Peeters, F.-X. Standaert, N. Donckers, J.-J. Quisquater, *Improved Higher-Order Side-Channel Attacks with FPGA Experiments*, CHES 2005, Lecture Notes in Computer Science, vol. 3659, pp. 309–323.
11. E. Prouff, *DPA Attacks and S-Boxes*, FSE 2005, LNCS, vol. 3557, pp. 424–441.
12. K. Schramm, C. Paar, *Higher Order Masking of the AES*, CT-RSA 2006, Lecture Notes in Computer Science, vol. 3860, pp. 208–225.
13. W. Schindler, K. Lemke, C. Paar, *A Stochastic Model for Differential Side-Channel Cryptanalysis*, CHES 2005, LNCS, vol. 3659, pp. 30–46.
14. F.-X. Standaert, T.G. Malkin, M. Yung, *A Unified Framework for the Analysis of Side-Channel Key Recovery Attacks*, Cryptology ePrint Archive, Report 2006/139, 2006, <http://eprint.iacr.org/>.

A Probability density function of the variable T_q

We take an iterative approach and first consider the intersection after two leakages. Assuming that the leakages respectively give rise to uniform lists of n_1 and n_2 candidates and the the key space has size $N = 2^n$, it yields $\Pr[T_2 = i | n_1, n_2] = \frac{\binom{n_1-1}{i-1} \cdot \binom{N-n_1}{n_2-i}}{\binom{N-1}{n_2-1}}$, where the binomials are taken among sets of $N-1$ possible elements since there is one fixed key that is not chosen uniformly. Then, assuming the knowledge of the distribution of $T_q(n_1, n_2, \dots, n_q)$ and an additional leakage that gives rise to a uniform list of n_{new} candidates, we can derive the distribution of T_{q+1} as follows: $\Pr[T_{q+1} = j | T_q = i, n_{new}] = \sum_i \Pr[T_{q+1} = j | T_q = i, n_{new}] \cdot \Pr[T_q = i]$, with: $\Pr[T_{q+1} = j | T_q = i, n_{new}] = \frac{\binom{i-1}{j-1} \cdot \binom{N-i}{n_{new}-j}}{\binom{N-1}{n_{new}-1}}$.

B Additional figures

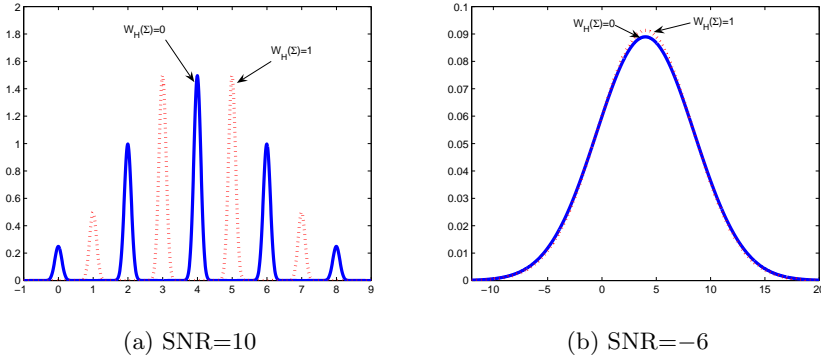


Fig. 13: Leakages PDFs in function of the noise: masked implementation.

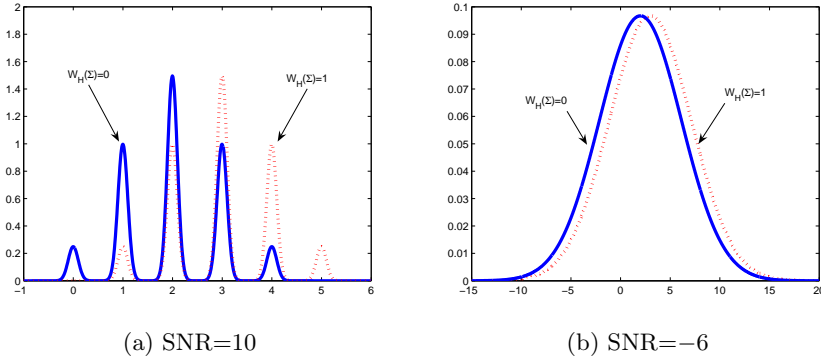


Fig. 14: Leakages PDFs in function of the noise: unmasked implementation.