

A Unified Framework for the Analysis of Side-Channel Key Recovery Attacks

François-Xavier Standaert^{1,*}, Tal G. Malkin², and Moti Yung^{2,3}

¹ UCL Crypto Group, Université catholique de Louvain

² Dept. of Computer Science, Columbia University

³ Google Inc.

fstandae@uclouvain.be, tal,moti@cs.columbia.edu

Abstract. The fair evaluation and comparison of side-channel attacks and countermeasures has been a long standing open question, limiting further developments in the field. Motivated by this challenge, this work makes a step in this direction and proposes a framework for the analysis of cryptographic implementations that includes a theoretical model and an application methodology. The model is based on commonly accepted hypotheses about side-channels that computations give rise to. It allows quantifying the effect of practically relevant leakage functions with a combination of information theoretic and security metrics, measuring the quality of an implementation and the strength of an adversary, respectively. From a theoretical point of view, we demonstrate formal connections between these metrics and discuss their intuitive meaning. From a practical point of view, the model implies a unified methodology for the analysis of side-channel key recovery attacks. The proposed solution allows getting rid of most of the subjective parameters that were limiting previous specialized and often ad hoc approaches in the evaluation of physically observable devices. It typically determines the extent to which basic (but practically essential) questions such as “*How to compare two implementations?*” or “*How to compare two side-channel adversaries?*” can be answered in a sound fashion.

1 Introduction

Traditionally, cryptographic algorithms provide security against an adversary who has only black box access to cryptographic devices. However, such a model does not always correspond to the realities of physical implementations. During the last decade, it has been demonstrated that targeting actual hardware rather than abstract algorithms may lead to very serious security issues. In this paper, we investigate the context of side-channel attacks, in which adversaries are enhanced with the possibility to exploit physical leakages such as power consumption [19] or electromagnetic radiation [2,14]. A large body of experimental work has been created on the subject and although numerous countermeasures are proposed in the literature, protecting implementations against such attacks

* Associate researcher of the Belgian Fund for Scientific Research (FNRS - F.R.S.)

is usually difficult and expensive. Moreover, most proposals we are aware of only increase the difficulty of performing the attacks, but do not fundamentally prevent them. Eventually, due to the device-specific nature of side-channel attacks, the comparison of their efficiency and the evaluation of leaking implementations are challenging issues, *e.g.* as mentioned in [22], page 163.

Following this state-of-the art, our work is mainly motivated by the need of having sound tools (*i.e.* a middle-ware between the abstract models and the concrete devices) to evaluate and compare different implementations and adversaries. As a matter of fact, the evaluation criteria in physically observable cryptography should be unified in the sense that they should be adequate and have the same meaning for analyzing any type of implementation or adversary. This is in clear contrast with the combination of ad hoc solutions relying on specific ideas designers have in mind. For example, present techniques for the analysis of side-channel attacks typically allow the statement of claims such as: “*An implementation X is better than an implementation Y against an adversary A*”. But such claims are of limited interest since an unsuccessful attack may theoretically be due both to the quality of the target device or to the ineffectiveness of the adversary. The results in this paper aim to discuss the extent to which more meaningful (adversary independent) statements can be claimed such as: “*An implementation X is better than an implementation Y*”. Similarly, when comparing different adversaries, present solutions for the analysis of side-channel attacks typically allow the statement of claims such as: “*An adversary A successfully recovers one key byte of an implementation X after the observation of q measurement queries.*”. But in practice, recovering a small set of key candidates including the correct one after a low number of measurement queries may be more critical for the security of an actual system than recovering the key itself after a high number of measurement queries (*e.g.* further isolating a key from a list can employ classical cryptanalysis techniques exploiting black box queries). The results in this paper aim at providing tools that help claiming more flexible statements and can capture various adversarial strategies.

Quite naturally, the previous goals imply the need of a sound model for the analysis of side-channel attacks. But perhaps surprisingly (and to the best of our knowledge), there have been only a few attempts to provably address physical security issues. A significant example is the work of Micali and Reyzin who initiated an analysis of side-channels taking the modularity of physically observable computations into account. The resulting model in [24] is very general, capturing almost any conceivable form of physical leakage. However and as observed by the authors themselves, this generality implies that the obtained positive results (*i.e.* leading to useful constructions) are quite restricted in nature and it is not clear how they apply to practice. This is especially true for primitives such as modern block ciphers for which even the black box security cannot be proven.

In the present work, we consequently give up a part of this generality and concentrate on current attacks (*i.e.* key recovery) and adversaries (*i.e.* statistical procedures to efficiently discriminate the key), trying to keep a sound and systematic approach aside these points. For this purpose, we first separate the implementation

issue (*i.e.* “*how good is my implementation?*”) and the adversarial issue (*i.e.* “*how strong is my adversary?*”) in the physically observable setting. We believe that the methodological division of both concerns brings essential insights and avoids previous confusions in the analysis of side-channel attacks. As a consequence, we introduce two different types of evaluation metrics. First, an information theoretic metric is used to measure the amount of information that is provided by a given implementation. Second, an actual security metric is used to measure how this information can be turned into a successful attack. We propose candidates for these metrics and show that they allow comparing different implementations and adversaries. We also demonstrate important connections between them in the practically meaningful context of Gaussian leakage distributions and discuss their intuitive meaning. Eventually, we move from formal definitions to practice-oriented definitions in order to introduce a unified evaluation methodology for side-channel key recovery attacks. We also provide an exemplary application of the model and discuss its limitations.

Related works include a large literature on side-channel issues, ranging from attacks to countermeasures and including statistical analysis concerns. The side-channel lounge [12], DPA book [22] and CHES workshops [8] provide a good list of references, a state-of-the art view of the field and some recent developments, respectively. Most of these previous results can be re-visited in the following framework in order to improve their understanding. The goal of this paper is therefore to facilitate the interface between theoretical and practical aspects in physically observable cryptography. We mention that in parallel to our work, the models in [3,20] consider a restricted context of noiseless leakages. They allow deriving formal bounds on the efficiency of certain attacks but are not aimed to analyze actual devices (that always have to deal with noise) which is our main goal. Finally, [25] initiated a study of forward secure cryptographic constructions with rigorous security analysis of side-channel attacks. [11,26] then proposed similar constructions in a more general setting and standard model. These works exploit assumptions such as bounded adversaries or leakages of which the validity can be measured for different devices thanks to the methodology in this paper.

Finally, our analysis employs ideas from the classical communication theory [10,28,29]. But whereas source and channel coding attempt to represent the information in an efficient format for transmission, cryptographic engineers have the opposite goal to make their circuit’s internal configurations unintelligible to the outside world. This analogy provides a rationale for our metrics. Note that different measures of uncertainty have frequently been used in the cryptographic literature to quantify the effectiveness of various attacks, *e.g.* in [6]. Our line of research follows a slightly different approach in the sense that we assign specific tasks to different metrics. Namely, we suggest to evaluate implementations with an information theoretic metric (conditional entropy) and to evaluate attacks and adversaries with security metrics (success rates or guessing entropy). This allows us to consider first implementations as non-adversarial information emitting objects where keys are randomly chosen, and then adversaries which operate under certain (computational and access) restrictions on top of the implementations. This

duality enables our model to be reflective of the situation in the real world and therefore to be useful beyond theoretical analysis, *i.e.* applicable to any simulated or actual lab data, for various cryptographic algorithms.

Note that because of place constraints, proofs and technical details have been removed from the paper and made available in an extended version [30].

2 Intuitive Description of the Model and Terminology

In this section, we give an intuitive description of side-channel key recovery attacks that will be formally defined and investigated in the rest of this paper.

A generic side-channel key recovery is illustrated in Figure 1 that we explain as follows. First, the term *primitive* is used to denote cryptographic routines corresponding to the practical instantiation of some idealized functions required to solve cryptographic problems. For example, the AES Rijndael is a cryptographic primitive. Second, the term *device* is used to denote the physical realization of a cryptographic primitive. For example, a smart card running the AES Rijndael can be the target device of a side-channel attack. A *side-channel* is an unintended communication channel that leaks some information from a device through a physical media. For example, the power consumption or the electromagnetic radiation of a target device can be used as side-channels. The output of a side-channel is a *physical observable*. Then, the *leakage function* is an abstraction that models all the specificities of the side-channel and the measurement setup used to monitor the physical observables. An *implementation* is the combination of a cryptographic device and a leakage function. Finally, a *side-channel adversary* is an algorithm (or a set of algorithms) that can query the implementation to get the leakage function results in addition to the traditional black-box access. Its goal is to defeat a given security notion (*e.g.* key recovery) within certain computational bounds and capabilities. Note that leakage functions and cryptographic implementations (aka physical computers) are formally defined in [24] and this paper relies on the same assumption as theirs.

Figure 1 suggests that, similarly to the classical communication theory, two aspects have to be considered (and quantified) in physically observable cryptography. First, actual implementations leak information, independently of the adversary exploiting it. The goal of our information theoretic metric is to measure the side-channel leakages in order to give a sound answer to the question: “*how to compare different implementations?*”. Second, an adversary analogous to a specific decoder exploits these leakages. The goal of our security metrics is to measure the extent to which this exploitation efficiently turns the information available into a key recovery. Security metrics are the counterpart of the Bit-Error-Rate in communication problems and aim to answer the question: “*how to compare different adversaries?*”. Interestingly, the figure highlights the difference between an actual adversary (of which the goal is simply to recover some secret data) and an evaluator (of which the goal is to analyze and understand the physical leakages). For example, comparing different implementations with an information theoretic metric is only of interest for an evaluator.

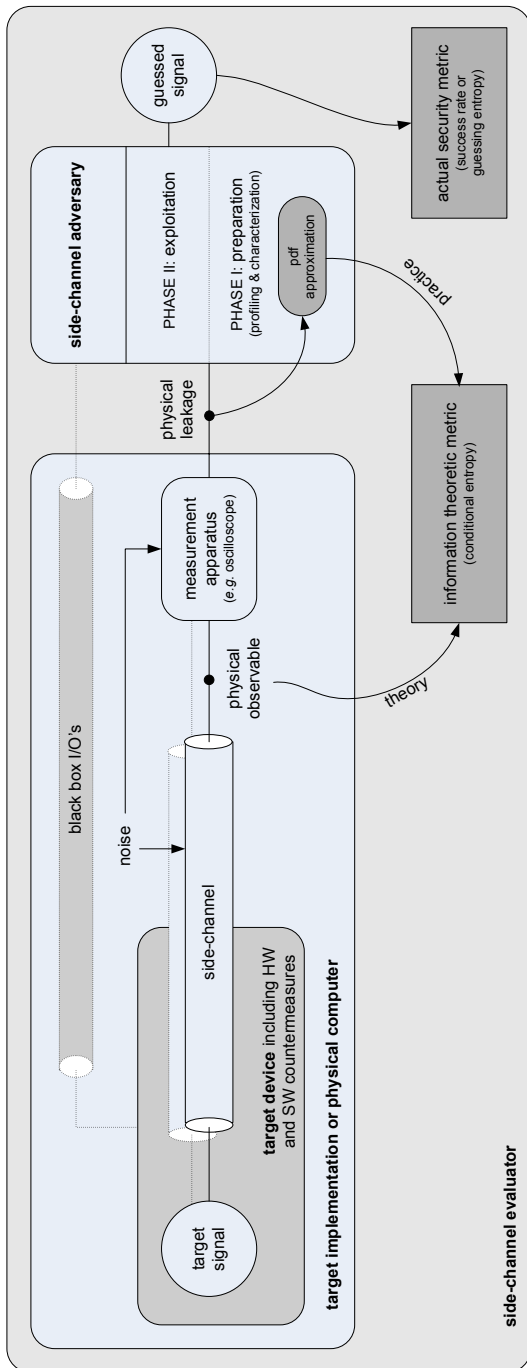


Fig. 1. Intuitive description of a side-channel key recovery attack

In practice, side-channel attacks are usually divided in two phases. First an (optional) preparation phase provides the adversary with a training device and allows him to profile and characterize the leakages. Second, an exploitation phase is directly mounted against the target device and is aimed to succeed the key recovery. Importantly, actual adversaries do not always have the opportunity to carry out a preparation phase (in which case profiling is done on the fly). By contrast, it is an important phase for evaluators since it allows performing optimized attacks and therefore leads to a better analysis of the physical leakages. Before moving to the definitions of our metrics, we finally mention the “theory” and “practice” arrows leading to the information theoretic metric in Figure 1. These arrows underline the fact that one can always assume a theoretical model for the side-channel and perform a *simulated attack*. If the model is meaningful, so is the simulated attack. But such simulations always have to be followed by an *experimental attack* in order to confirm the relevance of the model. Experimental attacks exploit actual leakages obtained from a measurement setup.

3 Formal Definitions

In this section, we define the metrics that we suggest for the analysis of physically observable devices. We first detail two possible security metrics, corresponding to different computational strategies. Both metrics relate to the notion of side-channel key recovery. Then, we propose an information theoretic metric driven by two requirements: (1) being independent of the adversary and (2) having the same meaning for any implementation or countermeasure. As a matter of fact and following the standard approach in information theory, Shannon’s conditional entropy is a good candidate for such a metric. Typically, the use of an average criteria to compare implementations is justified by the need of adversary independence. By contrast, the interactions of an adversary with a leaking system (*e.g.* adaptive strategies) are quantified with the security metrics in our model. We note that these candidate metrics will be justified by theoretical facts in Section 5 and practical applications in Section 6. However, it is an interesting open problem to determine if other metrics are necessary to evaluate side-channel attacks (*e.g.* min entropy is briefly discussed in Section 6).

3.1 Actual Security Metrics

Success Rate of the Adversary. Let $E_K = \{E_k(\cdot)\}_{k \in \mathcal{K}}$ be a family of cryptographic abstract computers indexed by a variable key K . Let (E_K, L) be the physical computers corresponding to the association of E_K with a leakage function L . As most cryptanalytic techniques, side-channel attacks are usually based on a divide-and-conquer strategy in which different (computationally tractable) parts of a secret key are recovered separately. In general, the attack defines a function $\gamma : \mathcal{K} \rightarrow \mathcal{S}$ which maps each key k onto an equivalent key class¹ $s = \gamma(k)$, such

¹ We focus on recovering key bytes for simplicity and because they are usual targets in side-channel attacks. But any other intermediate value in an implementation could be recovered, *i.e.* in general we can choose $s = \gamma(k, x)$ with x the input of $E_k(\cdot)$.

that $|\mathcal{S}| \ll |\mathcal{K}|$. We define a side-channel key recovery adversary as an algorithm $\mathbf{A}_{E_{K,L}}$ with time complexity τ , memory complexity m and q queries to the target physical computer. Its goal is to guess a key class $s = \gamma(k)$ with non negligible probability, by exploiting its collected (black box and physical) information. For this purpose, we assume that the output of the adversary $\mathbf{A}_{E_{K,L}}$ is a guess vector $\mathbf{g} = [g_1, g_2, \dots, g_{|\mathcal{S}|}]$ with the different key candidates sorted according to the attack result: the most likely candidate being g_1 . A practice-oriented description of $\mathbf{A}_{E_{K,L}}$ with a detailed specification of its features is given in [30], Appendix A. Finally, we define a side-channel key recovery of order o with the experiment:

Experiment $\mathbf{Exp}_{\mathbf{A}_{E_{K,L}}}^{\text{sc-kr-}o}$

```

 $\mathbf{g} \leftarrow \mathbf{A}_{E_{K,L}}; s = \gamma(k); k \xleftarrow{R} \mathcal{K};$ 
if  $s \in [g_1, \dots, g_o]$  then return 1;
else return 0;
```

The o^{th} -order success rate of $\mathbf{A}_{E_{K,L}}$ against a key class variable S is defined as:

$$\mathbf{Succ}_{\mathbf{A}_{E_{K,L}}}^{\text{sc-kr-}o,S}(\tau, m, q) = \Pr [\mathbf{Exp}_{\mathbf{A}_{E_{K,L}}}^{\text{sc-kr-}o} = 1] \quad (1)$$

Intuitively, a success rate of order 1 (*resp.* 2, ...) relates to the probability that the correct key is sorted first (*resp.* among the two first ones, ...) by the adversary. When not specified, a first order success rate is assumed.

Computational Restrictions. Similarly to black box security, computational restrictions have to be imposed to side-channel adversaries in order to capture the reality of physically observable cryptographic devices. This is the reason for the parameters τ, m, q . Namely, the attack time complexity τ and memory complexity m (mainly dependent on the number of key classes $|\mathcal{S}|$) are limited by present computer technologies. The number of measurement queries q is limited by the adversary's ability to monitor the device. In practice, these quantities are generally separated for the preparation and exploitation phases (see Section 5). But additionally to the computational cost of the side-channel attack itself, another important parameter is the remaining workload after the attack. For example, considering a success rate of order o implies that the adversary still has a maximum of o key candidates to test after the attack. If this has to be repeated for different parts of the key, it may become a non negligible task. As a matter of fact, the previously defined success rate measures an adversary with a fixed maximum workload after the side-channel attack. A more flexible metric that is also convenient in our context is the guessing entropy. It measures the average number of key candidates to test after the side-channel attack. The guessing entropy was originally defined in [23] and has been proposed to quantify the effectiveness of adaptive side-channel attacks in [20]. It can be related to the notion of gain that has been used in the context of multiple linear cryptanalysis to measure how much the complexity of an exhaustive key search is reduced thanks to an attack [5]. We use it as an alternative to the success rate.

Guessing Entropy. We first define a side-channel key guessing experiment:

$$\begin{aligned} &\text{Experiment } \mathbf{Exp}_{\mathbf{A}_{E_K, L}}^{\text{sc-kg}} \\ &[\mathbf{g} \leftarrow \mathbf{A}_{E_K, L}; s = \gamma(k); k \xleftarrow{R} \mathcal{K};] \\ &\text{return } i \text{ such that } g_i = s; \end{aligned}$$

The guessing entropy of $\mathbf{A}_{E_K, L}$ against a key class variable S is then defined as:

$$\mathbf{GE}_{\mathbf{A}_{E_K, L}}^{\text{sc-kr-}S}(\tau, m, q) = \mathbf{E}(\mathbf{Exp}_{\mathbf{A}_{E_K, L}}^{\text{sc-kg}}) \tag{2}$$

3.2 Information Theoretic Metric

Let S be the previously used target key class discrete variable of a side-channel attack and s be a realization of this variable. Let $\mathbf{X}_q = [X_1, X_2, \dots, X_q]$ be a vector of variables containing a sequence of inputs to the target physical computer and $\mathbf{x}_q = [x_1, x_2, \dots, x_q]$ be a realization of this vector. Let \mathbf{L}_q be a random vector denoting the side-channel observations generated with q queries to the target physical computer 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 L corresponding to the input vector \mathbf{x}_q . Let finally $\Pr[s|\mathbf{l}_q]$ be the conditional probability of a key class s given a leakage \mathbf{l}_q . We define the conditional entropy matrix as:

$$\mathbf{H}_{s, s^*}^q = - \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s^*|\mathbf{l}_q], \tag{3}$$

where s and s^* respectively denote the correct key class and a candidate out of the $|\mathcal{S}|$ possible ones. From 3, we derive Shannon’s conditional entropy:

$$\mathbf{H}[S|\mathbf{L}_q] = - \sum_s \Pr[s] \sum_{\mathbf{l}_q} \Pr[\mathbf{l}_q|s] \cdot \log_2 \Pr[s|\mathbf{l}_q] = \mathbf{E}_s(\mathbf{H}_{s, s}^q) \tag{4}$$

It directly yields the mutual information: $\mathbf{I}(S; \mathbf{L}_q) = \mathbf{H}[S] - \mathbf{H}[S|\mathbf{L}_q]$. Note that the inputs and outputs of an abstract computer are generally given to the side-channel adversary (but hidden in the formulas for clarity reasons). Therefore, it is implicitly a computational type of entropy that is proposed to evaluate the physical leakages. This is because divide-and-conquer strategies target a key class assuming that the rest of the key is unknown. But from a purely information theoretic point of view, the knowledge of a plaintext-ciphertext pair can determine a key completely (*e.g.* for block ciphers). Hence and as detailed in the next section, the amount of information extracted by a side-channel adversary depends on its computational complexity. Note also that leakage functions can be discrete or (most frequently) continuous. In the latter case, it is formally a conditional differential entropy that is computed. Note finally that in simulated attacks where an analytical model for a continuous leakage probability distribution is assumed, the previous sums over the leakages can be turned into integrals.

4 Practical Limitations

One important goal of the present framework is to allow a sound evaluation of any given implementation, if possible independently of an adversary’s algorithmic details. For this purpose, the strategy we follow is to consider an information theoretic metric that directly depends on the leakages probability distribution $\Pr[\mathbf{L}_q|S]$. Unfortunately, there are two practical caveats in this strategy.

First, the conditional probability distribution $\Pr[\mathbf{L}_q|S]$ is generally unknown. It can only be approximated through physical observations. This is the reason for the leakage function abstraction in the model of Micali and Reyzin that we follow in this work. It informally states that the only way an adversary knows the physical observables is through measurements. Therefore, practical attacks and evaluations have to exploit an approximated distribution $\hat{\Pr}[\mathbf{L}_q|S]$ rather than the actual one $\Pr[\mathbf{L}_q|S]$. Second, actual leakages may have very large dimensions since they are typically the output of a high sampling rate acquisition device like an oscilloscope. As a consequence, the approximation of the probability distribution for all the leakage samples is computationally intensive. Practical attacks usually approximate the probability distribution of a reduced set of samples, namely $\hat{\Pr}[\tilde{\mathbf{L}}_q|S]$. We denote side-channel attacks that exploit the approximated probability distribution of a reduced set of leakage samples as generic template attacks. A straightforward consequence of the previous practical limitations is that for any actual device, the mutual information $I(S; \mathbf{L}_q)$ can only be approximated through statistical sampling, by using generic template attacks.

We note there are different concerns in the application of template attacks such as: “how to limit the number of leakage samples for which the distribution will be estimated?” or “how to limit the number of templates to build?”. The data dimensionality reduction techniques used in [4,32] and the stochastic models in [16,27] can be used to answer these questions in a systematic manner. But there is no general theory allowing one to decide what is the best attack for a given device. Hence, in the following we will essentially assume that one uses the “best available tool” to approximate the leakage distribution. Quite naturally, the better generic template attacks perform in practice, the better our framework allows analyzing physical information leakages.

5 Relations between the Evaluation Metrics

In this section, we provide theoretical arguments that justify and connect the previous information theoretic and security metrics. These connections allow us to put forward interesting features and theoretical limitations of our model. In particular, we will consider two important questions.

First, as mentioned in Section 4, generic template attacks require to estimate the leakage probability distribution. Such a leakage model is generally built during a preparation phase and then used to perform a key recovery during an exploitation phase (as pictured in Figure 1). And as mentioned in Section 3.1, these phases have to be performed within certain computational limits. Hence,

to the previously defined complexity values τ, m, q of the online phase, one has to add the complexities of the preparation phase, denoted as τ_p, m_p, q_p . The first question we tackle is: given some bounds on (τ_p, m_p, q_p) , can an adversary build a good estimation of the leakage distribution? We show in Section 5.1 that the conditional entropy matrix of Equation (3) is a good tool to answer this question. We also show how it relates to the asymptotic success rate of a Bayesian adversary. Then, assuming that one can build a good approximation for the leakage distribution, we investigate the extent to which the resulting estimation of the mutual information allows comparing different implementations. Otherwise said, we analyze the dependencies between our information theoretic and security metrics. We show that there exist practically meaningful contexts of Gaussian side-channels for which strong dependencies can be put forward. But we also emphasize that no general statements can be made for arbitrary distributions. Section 5.2 essentially states that the mutual information is a good metric to compare different implementations, but it always has to be completed with a security analysis (*i.e.* success rate and/or guessing entropy).

5.1 Asymptotic Meaning of the Conditional Entropy: “Can I Approximate the Leakage Probability Distribution?”

We start with three definitions.

Definition 1. The asymptotic success rate of a side-channel adversary $\mathbf{A}_{\mathbf{E}_{K,L}}$ against a key class variable S is its success rate when the number of measurement queries q tends to the infinity. It is denoted as: $\mathbf{Succ}_{\mathbf{A}_{\mathbf{E}_{K,L}}}^{\text{sc-kr-}o,S}(q \rightarrow \infty)$.

Definition 2. Given a leakage probability distribution $\Pr[\mathbf{L}_q|S]$ and a number of side-channel queries stored in a leakage vector \mathbf{l}_q , a Bayesian side-channel adversary is an adversary that selects the key as $\mathit{argmax}_{s^*} \Pr[s^*|\mathbf{l}_q]$.

Definition 3. An approximated leakage distribution $\tilde{\Pr}[\tilde{\mathbf{L}}_q|S]$ is sound if the first-order asymptotic success rate of a Bayesian side-channel adversary exploiting this leakage distribution against the key class variable S equals one.

In this section, we assume that one has built an approximated leakage distribution $\tilde{\Pr}[\tilde{\mathbf{L}}_q|S]$ with some (bounded) measurement queries q_p , memory m_p and time τ_p . We want to evaluate if this approximation is good. For theoretical purposes, we consider an adversary/evaluator who can perform unbounded queries to the target device during the exploitation phase. We use these queries to evaluate the entropy matrix $\hat{\mathbf{H}}_{s,s^*}^q$ defined in Section 3.2. It directly leads to the following relation with the asymptotic success rate of a Bayesian adversary.

Theorem 1. *Assuming independent leakages for the different queries in a side-channel attack, an approximated leakage probability distribution $\tilde{\Pr}[\tilde{\mathbf{L}}_q|S]$ is sound if and only if the conditional entropy matrix evaluated in an unbounded exploitation phase is such that $\mathit{argmin}_{s^*} \hat{\mathbf{H}}_{s,s^*}^q = s, \forall s \in \mathcal{S}$.*

The proof of Theorem 1 is given in [30]. There are several important remarks:

1. Theorem 1 only makes sense for bounded preparation phases. For unbounded preparations, an adversary would eventually access the exact distribution $\Pr[\mathbf{L}_q|S]$. In this context, the soundness does only depend on the cardinality of the different sets $\{s^* | \Pr[\mathbf{L}_q|s^*] = \Pr[\mathbf{L}_q|s]\}$, $\forall s \in \mathcal{S}$.
2. The condition of independence for consecutive leakages is not expected to be fully verified in practice. For example, there could exist history effects in the side-channel observations. However, it is expected to hold to a sufficient degree for our proof to remain meaningful in most applications.
3. In practice, the exploitation phase in a side-channel attack is bounded as the preparation. Therefore, Theorem 1 will be relevant as long as the number of leakages used to test the approximated leakage distribution and estimate the conditional entropy matrix is sufficient.
4. Finally, the condition on the entropy matrix $\widehat{\mathbf{H}}_{s,s^*}^q$ is stated for the number of queries q for which the leakage distribution $\Pr[\mathbf{L}_q|S]$ was approximated during the preparation phase. In general, finding a sound approximation for q implies that it should also be feasible to find sound approximations for any $q' > q$. But in practice, computational limitations can make it easier to build a sound approximation for small q values than for larger ones.

5.2 Comparative Meaning of the Conditional Entropy: “Does More Entropy Imply More Security?”

Let us write an exemplary conditional entropy matrix and its estimation as:

$$\mathbf{H}_{s,s^*}^q = \begin{pmatrix} h_{1,1} & h_{1,2} & \dots & h_{1,|S|} \\ h_{2,1} & h_{2,2} & \dots & h_{2,|S|} \\ \dots & \dots & \dots & \dots \\ h_{|S|,1} & h_{|S|,2} & \dots & h_{|S|,|S|} \end{pmatrix} \quad \widehat{\mathbf{H}}_{s,s^*}^q = \begin{pmatrix} \hat{h}_{1,1} & \hat{h}_{1,2} & \dots & \hat{h}_{1,|S|} \\ \hat{h}_{2,2} & \hat{h}_{2,2} & \dots & \hat{h}_{2,|S|} \\ \dots & \dots & \dots & \dots \\ \hat{h}_{|S|,1} & \hat{h}_{|S|,2} & \dots & \hat{h}_{|S|,|S|} \end{pmatrix}$$

Theorem 1 states that if the diagonal values of a (properly approximated) matrix are minimum for all key classes $s \in \mathcal{S}$, then these key classes can be asymptotically recovered by a Bayesian adversary. As a matter of fact, it gives rise to a binary conclusion about the approximated leakage probability distribution. Namely, Theorem 1 answers the question: “*Can one approximate the leakage probability distribution under some computational bounds τ_p, m_p, q_p ?*”

Let us now assume that the answer is positive and denote each element $h_{s,s}$ as the residual entropy of a key class s . In this subsection, we are interested in the values of these entropy matrix elements. In particular, we aim to highlight the relation between these values and the effectiveness of a side-channel attack, measured with the success rate. Otherwise said, we are interested in the question: “*Does less entropy systematically implies a faster convergence towards a 100% success rate?*”. Contrary to the previous section, this question makes sense both for the ideal conditional entropy matrix that would correspond to an exact leakage distribution and for its approximation. Since general conclusions for arbitrary leakage distributions are not possible to obtain, our strategy is to

first consider simple Gaussian distributions and to extrapolate the resulting conclusions towards more complex cases. We start with three definitions.

Definition 4. An $|\mathcal{S}|$ -target side-channel attack is an attack where an adversary tries to identify one key class s out of $|\mathcal{S}|$ possible candidates.

Definition 5. An univariate (*resp.* multivariate) leakage distribution is a probability distribution predicting the behavior of one (*resp.* several) leakage samples.

Definition 6. A Gaussian leakage distribution is the probability distribution of a leakage function that can be written as the sum of a deterministic part and a normally distributed random part, with mean zero and standard deviation σ .

Finally, since we now consider the residual entropies of the different key classes separately, we need a more specific definition of the success rate against a key class s (*i.e.* a realization of the variable S), denoted as $\text{Succ}_{\mathcal{A}_{E_{k,l}}}^{\text{sc-kr-}o,s}(\tau, m, q)$. It corresponds to the definition of Section 3.1 with a fixed key class.

Examples. Figure 2 illustrates several Gaussian leakage distributions. The upper left picture represents the univariate leakage distributions of a 2-target side-channel attack, each Gaussian curve corresponding to one key class s . The upper right picture represents the bivariate leakage distributions of a 2-target side-channel attack. Finally, the lower left and right pictures represent the univariate and bivariate leakage distributions of an 8-target side-channel attack. Note that in general, the internal state of an implementation does not only depend on the keys but also on other inputs, *e.g.* the plaintexts in block ciphers. Hence, the different dimensions in a multivariate distribution can represent both the different samples of a single leakage trace (generated with a single plaintext) or different traces (*e.g.* each dimension could correspond to a different plaintext). Eventually, it is an adversary's choice to select the internal states for which templates will be built. Therefore, we do not claim that these distributions always connect to practical attacks. But as will be seen in the following, even these simple theoretical contexts hardly allow simple connections between information and security.

We now discuss formally the connections between the success rate against a key class s and its residual entropy for idealized distributions and attacks.

Definition 7. An ideal side-channel attack is a Bayesian attack in which the leakages are exactly predicted by the adversary's approximated probability density function $\hat{\text{Pr}}[\tilde{\mathbf{L}}_q|S]$ (*e.g.* thanks to an unbounded preparation phase).

Lemma 1. *In an ideal 2-target side-channel attack exploiting a univariate Gaussian leakage distribution, the residual entropy of a key class s is a monotonously decreasing function of the single query (hence multi-queries) success rate against s .*

Lemma 2. *In an ideal 2-target side-channel attack exploiting a multivariate Gaussian leakage distribution, with independent leakage samples having the same noise standard deviation, the residual entropy of a key class s is a monotonously decreasing function of the single query (hence multi-queries) success rate against s .*

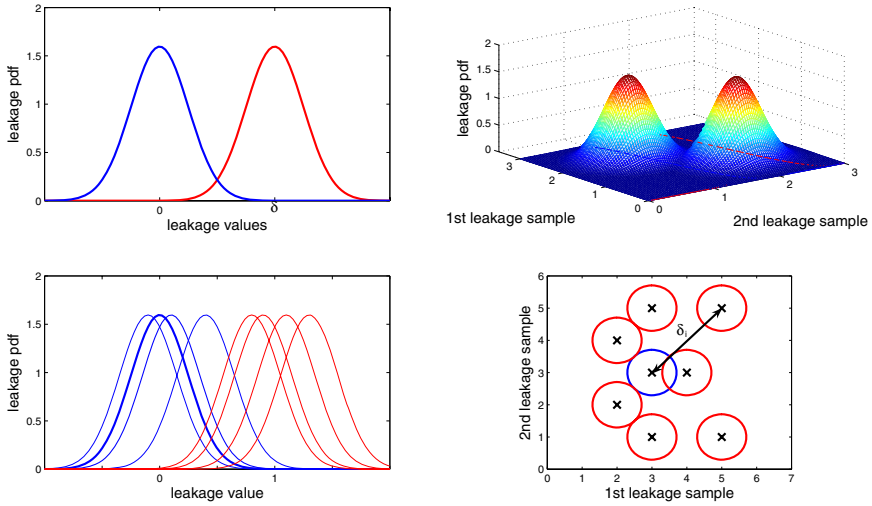


Fig. 2. Illustrative leakage probability distributions $\Pr[\mathbf{L}_q|S]$

These lemmas essentially state that (under certain conditions) the entropy and success rate in a 2-target side-channel attack only depend on the distance between the target leakages mean values normalized by their variance. It implies the direct intuition that more entropy means less success rate. Unfortunately, when moving to the $|\mathcal{S}|$ -target case with $|\mathcal{S}| > 2$, such a perfect dependency does not exist anymore. One can observe in the lower right part of Figure 2 that the entropy and success rate not only depend on the normalized distances δ_i/σ but also on how the keys are distributed within the leakage space. Therefore, we now define a more specific context in which formal statements can be proven.

Definition 8. A perfect Gaussian leakage distribution $\Pr[\mathbf{L}_q|s]$ for a key class s is a Gaussian leakage distribution with independent leakage samples having the same noise standard deviation such that the Euclidean distance between each key class candidate mean value and the correct key class candidate mean value is equal and the residual entropy of the key class s is maximum.

Theorem 2. *In an ideal side-channel attack exploiting a perfect Gaussian leakage distribution, the residual entropy of a key class s is a monotonously decreasing function of the single query (hence multi-queries) success rate against s .*

The proofs of Lemmas 1, 2 and Theorem 2 are given in [30]. They constitute our main positive results for the use of the conditional entropy as a comparison metric for different implementations. Unfortunately, in the most general context of non perfect leakage distributions, those general statements do not hold. Facts 1 and 2 in [30] even demonstrate that there exist no generally true dependencies between the conditional entropy and the success rate in a general setting.

5.3 Intuition of the Metrics

In this section, we recall and detail a number of important intuitions that can be extracted from the previous theory. We also discuss how they can be exploited in practical applications and highlight their limitations.

Intuitions Related to Theorem 1

- 1.1 *Theorem 1 tells if it is possible to approximate a given leakage function in a bounded preparation phase.* As mentioned in Section 4, such an approximation highly depends on the actual tools that are used for this purpose. In general, the better the tools, the better the evaluation. Hence, Theorem 1 allows checking if these tools are powerful enough. If they are not...
- 1.2 *Theorem 1 indicates some resistance of the target implementation against side-channel attacks.* If one cannot build a sound approximation of the leakage probability distribution, even with intensive efforts, then the 1st-order asymptotic success rate of the Bayesian side-channel adversary does not reach one. But this does not imply security against side-channel attacks (*e.g.* think about a device where only one key could not be recovered). In this context, it is important to evaluate the actual security metrics for different adversaries in order to check if high success rates can still be reached.

Intuitions Related to Theorem 2

- 2.1 *Theorem 2 only applies to sound leakage distributions.* Intuitively, it means that comparing the conditional entropy provided by different leakage functions only makes sense if the corresponding approximated leakage probability distribution lead to asymptotically successful attacks.
- 2.2 *Theorem 2 confirms that mutual information is a relevant tool to compare different implementations.* It shows meaningful contexts of Gaussian channels for which less residual entropy for a key class implies a more efficient attack. It strengthens the intuitive requirements of Section 3, namely an adversary independent metric with the same meaning for any implementation.
- 2.3 *The conditional entropy is not a stand-alone metric to compare implementations and always has to be combined with a security analysis.* This relates both to theoretical limitations (since there exists no general relation between information and security) and practical constraints. For a given amount of information leaked by an implementation, different side-channel distinguishers could be considered. Therefore, security metrics are useful to evaluate the number of queries for these different attacks to succeed.

Note that the mutual information, success rates and guessing entropy are average evaluation criteria. However in practice, the information leakages and security of an implementation could be different for different keys. Therefore, it is important to also consider these notions for the different keys separately (*e.g.* to evaluate

the conditional entropy matrix rather than the mutual information). This last remark motivates the following practice-oriented definition.

Definition 9. We say that a side-channel attack against a key class variable S is a weak template attack if all the key classes s have the same residual entropy $h_{s,s}$ and each line of the entropy matrix $\mathbf{H}_{s,s}^q$ is a permutation of another line of the matrix. We say that a side-channel attack is a strong template attack if at least one of the previous conditions does not hold.

6 Applications of the Model

In order to confirm that (although limited by theoretical concerns) the intuition of Theorem 2 applies to practice, this section provides examples of side-channel attacks that can be reproduced by the reader. Applications to more complex and practically meaningful contexts can be found in other publications [21,30,31,32,33].

For this purpose, we consider a known plaintext attack against a reduced block cipher that we formalize as follows. Let \mathbf{S} be a 4-bit substitution box, *e.g.* from the AES candidate Serpent. We target the computation of $y = \mathbf{S}(x \oplus k)$, where x is a random plaintext and k a secret key. A Bayesian adversary is provided with observations $(x, L'(y) + r)$, where r is a gaussian noise with mean 0 and standard deviation σ . For any y value, the deterministic part of the leakage $L'(y)$ is given by a vector \mathbf{Z} . The adversary's goal is to recover the key k . Our simulations exploit different leakage functions and assume an unbounded preparation phase (*i.e.* the adversary has knowledge of the exact leakage distribution). We start the frequently observed Hamming weight leakages and $\mathbf{Z}_1 = [0, 1, 1, 2, 1, 2, 2, 3, 1, 2, 2, 3, 2, 3, 3, 4]$. We also evaluate two other leakage functions represented by the vectors: $\mathbf{Z}_2 = [0, 0, 0, 0, 1, 1, 1, 1, 2, 2, 2, 2, 3, 3, 3, 3]$ and $\mathbf{Z}_3 = [0, 0, 0, 0, 0, 0, 0, 0, 1, 2, 2, 3, 3, 4, 4]$. The conditional entropies and single-query success rates with $\sigma = 0$ can be straightforwardly computed as:

$$\begin{aligned}
 H[K|\mathbf{L}_1^{\mathbf{Z}_1}] &\simeq 1.97 & H[K|\mathbf{L}_1^{\mathbf{Z}_2}] &= 2 & H[K|\mathbf{L}_1^{\mathbf{Z}_3}] &\simeq 2.16 \\
 \text{Succ}_{\mathbf{L}_1^{\mathbf{Z}_1}}^{\text{sc-kr}}(q=1) &= \frac{5}{16} & \text{Succ}_{\mathbf{L}_1^{\mathbf{Z}_2}}^{\text{sc-kr}}(q=1) &= \frac{1}{4} & \text{Succ}_{\mathbf{L}_1^{\mathbf{Z}_3}}^{\text{sc-kr}}(q=1) &= \frac{5}{16}
 \end{aligned}$$

At first sight, it seems that these leakage functions exactly contradict Theorem 2. For example, when moving from \mathbf{Z}_2 to \mathbf{Z}_3 , we see that both the conditional entropy and the success rate are increased. However, the goal of side-channel attacks is generally to reach high success rates that are not obtained with a single query. Hence, it is also interesting to investigate the success rate for more queries. In the left part of Figure 3, these success rates for increasing q values are plotted. It clearly illustrates that while \mathbf{Z}_2 leads to a lower success rate than \mathbf{Z}_3 for $q = 1$, the opposite conclusion holds when increasing q . That is, the intuition given by Theorem 2 only reveals itself for $q > 2$. Importantly, these conclusions can vary when noise is inserted in the leakages, *e.g.* assuming $\sigma = 1$, we have:

$$\begin{aligned}
 H[K|\mathbf{L}_1^{\mathbf{Z}_1}] &\simeq 3.50 & H[K|\mathbf{L}_1^{\mathbf{Z}_2}] &\simeq 3.42 & H[K|\mathbf{L}_1^{\mathbf{Z}_3}] &\simeq 3.22
 \end{aligned}$$

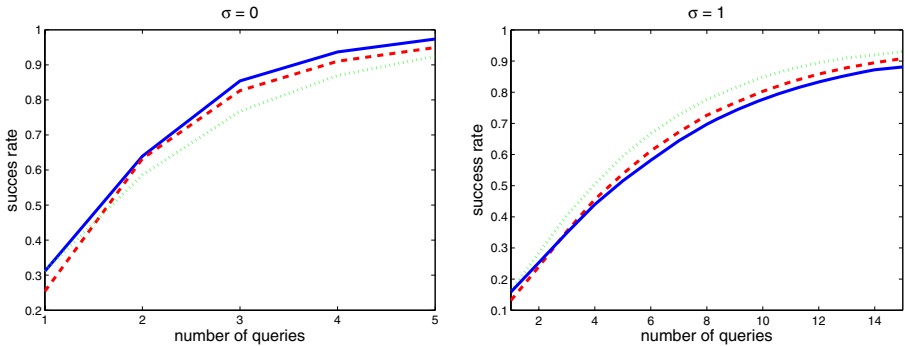


Fig. 3. 1st-Order success rates in function of the number of queries for the leakages functions corresponding to \mathbf{Z}_1 (solid line), \mathbf{Z}_2 (dashed line) and \mathbf{Z}_3 (dotted line)

The right part of Figure 3 plots the success rates of these noisy leakage functions. It again highlights a context in which Theorem 2 is eventually respected. In general, these examples underline another important feature of our metrics. Namely, the more challenging the side-channel attack (*i.e.* the more queries needed to reach high success rates), the more significant the conditional entropy is. Otherwise said: the mutual information better reveals its intuition asymptotically. And in such contexts, the single-query success rate can be misleading.

Note that the examples in this section are more directly reflective of actual side-channel attacks in which different plaintexts can generally be used to identify a key class than the ideal contexts investigated in Section 5.2.

A Short Note on Minimum Entropy. With respect to the relevance of other metrics in the model, we finally mention that min entropy is equivalent to a single-query success rate. Since side-channel attacks are essentially multiple-query attacks, we believe that Shannon’s conditional entropy better captures the information leakages in most practical applications. For example, Figure 3 is typical of contexts where min entropy is misleading, *i.e.* where the success rate for $q = 1$ is not very significant while the conditional entropy nicely quantifies the evolution of this success rate for any larger q . But as already said, the information theoretic analysis always has to be completed with a security analysis. Hence, even in contexts where min entropy is the right metric, our model would detect it.

7 Evaluation Methodology

Following the previous sections, an evaluation methodology for side-channel attacks intends to analyze both the quality of an implementation and the strength of an adversary, involving the five steps illustrated in Figure 4. It again indicates that the information theoretic metric can be used to measure an implementation while the actual security metrics are rather useful to evaluate adversaries. Additionally to these metrics, it is often interesting to define a Signal-to-Noise Ratio (SNR) in order to determine the amount of noise in the physical observations.

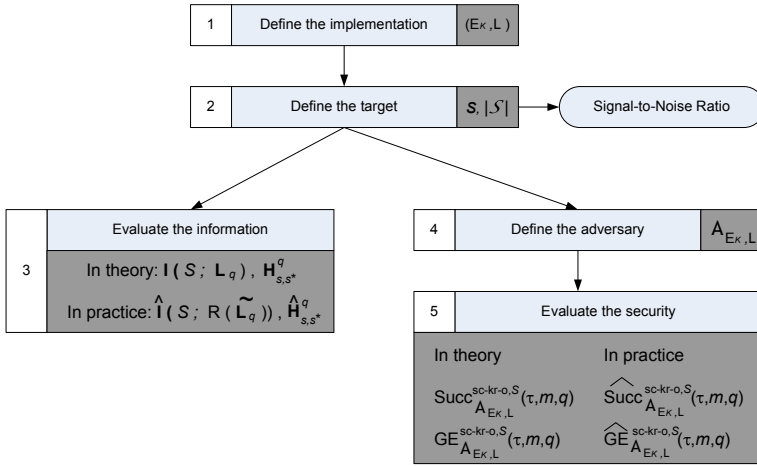


Fig. 4. Evaluation methodology for side-channel attacks

Since noise insertion is a generic countermeasure to improve resistance against side-channel attacks, it can be used to plot the information theoretic and security metrics with its respect. We note finally that the definition of an implementation requires to evaluate the cost of the equipment used to monitor the leakages. Since quantifying such costs is typically the tasks assigned the standardization bodies, we refer to the common criteria [9] and FIPS 140-2 documents [13] (or alternatively to the IBM taxonomy [1]) for these purposes. In general, the benefit of the present model is not to solve these practical issues but to state the side-channel problem in a sound framework for its analysis. Namely, it is expected that the proposed security and information theoretic metrics can be used for the fair analysis, evaluation and comparison of any physical implementation or countermeasure against any type of side-channel attack.

8 Conclusions and Open Problems

A framework for the analysis of cryptographic implementations is introduced in order to unify the theory and practice of side-channel attacks. It is aimed to bridge the formal understanding of physically observable cryptography to the exploitation of actual leakages in experimental key recoveries. The framework is centered around a theoretical model in which the effect of practically relevant leakage functions is evaluated with a combination of security and information theoretic metrics. It allows discussing the underlying tradeoffs in physically observable cryptography in a fair manner. As an interface between an engineering problem (how much is leaked?) and a cryptographic problem (how to exploit it?), our framework helps putting forward properly quantified weaknesses in physically observable devices. The fair evaluations that it provides can then be used in two directions. Either the physical weaknesses can be sent to hardware designers

in order to reduce physical leakages. Or they can be transmitted to cryptographic designers in order to conceive schemes that can cope with physical leakages.

Open questions derive from this model in different directions. A first one relates to the best exploitation of large side-channel traces, *i.e.* to the construction of (ideally) optimal distinguishers. This requires investigating the best heuristics to deal with high dimensional leakage data (our model assumes adversaries exploiting such specialized algorithms). A second one relates to the investigation of stronger security notions than side-channel key recovery. That is, the different security notions considered in the black box model (*e.g.* undistinguishability from an idealized primitive) should be considered in the physical world, as initiated in [24] (but again in a more specialized way). A third possible direction relates to the construction of implementations with provable (or arguable) security against side-channel attacks, *e.g.* as proposed in [11,26,25]. Finally, this work could be extended to other physical threats (*e.g.* fault attacks) and combined with other approaches for modeling physical attacks such as [15,17,18].

References

1. Abraham, D.G., Dolan, G.M., Double, G.P., Stevens, J.V.: Transaction Security System. IBM Systems Journal 30(2), 206–229 (1991)
2. Agrawal, D., Archambeault, B., Rao, J., Rohatgi, P.: The EM side-channel(s). In: Kaliski Jr., B.S., Koç, Ç.K., Paar, C. (eds.) CHES 2002. LNCS, vol. 2523, pp. 29–45. Springer, Heidelberg (2003)
3. Backes, M., Köpf, B.: Formally Bounding the Side-Channel Leakage in Unknown-Message Attacks, IACR ePrint archive (2008), <http://eprint.iacr.org/2008/162>
4. Archambeau, C., Peeters, E., Standaert, F.-X., Quisquater, J.-J.: Template attacks in principal subspaces. In: Goubin, L., Matsui, M. (eds.) CHES 2006. LNCS, vol. 4249, pp. 1–14. Springer, Heidelberg (2006)
5. Biryukov, A., De Cannière, C., Quisquater, M.: On multiple linear approximations. In: Franklin, M. (ed.) CRYPTO 2004. LNCS, vol. 3152, pp. 1–22. Springer, Heidelberg (2004)
6. Cachin, C.: Entropy Measures and Unconditional Security in Cryptography, PhD Thesis, ETH Dissertation, num 12187, Zurich, Switzerland (1997)
7. Chari, S., Rao, J., Rohatgi, P.: Template attacks. In: Kaliski Jr., B.S., Koç, Ç.K., Paar, C. (eds.) CHES 2002. LNCS, vol. 2523, pp. 13–28. Springer, Heidelberg (2003)
8. Cryptographic Hardware and Embedded Systems, <http://www.chesworkshop.org>
9. Application of Attack Potential to Smart Cards, Common Criteria Supporting Document, Version 1.1 (July 2002), <http://www.commoncriteriaportal.org>
10. Cover, T.M., Thomas, J.A.: Information Theory. Wiley and Sons, New York (1991)
11. Dziembowski, S., Pietrzak, K.: Leakage-Resilient Cryptography. In: The proceedings of FOCS 2008, Philadelphia, USA, pp. 293–302 (October 2008)
12. ECRYPT Network of Excellence in Cryptology, The Side-Channel Cryptanalysis Lounge, http://www.crypto.ruhr-uni-bochum.de/en_sclounge.html
13. FIPS 140-2, Security Requirements for Cryptographic Modules, Federal Information Processing Standard, NIST, U.S. Dept. of Commerce (December 3, 2002)
14. Gandolfi, K., Mourtel, C., Olivier, F.: Electromagnetic analysis: Concrete results. In: Koç, Ç.K., Naccache, D., Paar, C. (eds.) CHES 2001. LNCS, vol. 2162, pp. 251–261. Springer, Heidelberg (2001)

15. Gennaro, R., Lysyanskaya, A., Malkin, T.G., Micali, S., Rabin, T.: Algorithmic Tamper-Proof Security: Theoretical Foundations for Security Against Tampering. In: Naor, M. (ed.) TCC 2004. LNCS, vol. 2951, pp. 258–277. Springer, Heidelberg (2004)
16. Gierlichs, B., Lemke-Rust, K., Paar, C.: Templates vs. Stochastic methods. In: Goubin, L., Matsui, M. (eds.) CHES 2006. LNCS, vol. 4249, pp. 15–29. Springer, Heidelberg (2006)
17. Ishai, Y., Sahai, A., Wagner, D.: Private circuits: Securing hardware against probing attacks. In: Boneh, D. (ed.) CRYPTO 2003. LNCS, vol. 2729, pp. 463–481. Springer, Heidelberg (2003)
18. Ishai, Y., Prabhakaran, M., Sahai, A., Wagner, D.: Private circuits II: Keeping secrets in tamperable circuits. In: Vaudenay, S. (ed.) EUROCRYPT 2006. LNCS, vol. 4004, pp. 308–327. Springer, Heidelberg (2006)
19. Kocher, P.C., Jaffe, J., Jun, B.: Differential power analysis. In: Wiener, M. (ed.) CRYPTO 1999. LNCS, vol. 1666, pp. 398–412. Springer, Heidelberg (1999)
20. Köpf, B., Basin, D.: an Information Theoretic Model for Adaptive Side-Channel Attacks. In: The proceedings of ACMCCS 2007, Alexandria, VA, USA (October 2007)
21. Macé, F., Standaert, F.-X., Quisquater, J.-J.: Information theoretic evaluation of side-channel resistant logic styles. In: Paillier, P., Verbauwhede, I. (eds.) CHES 2007. LNCS, vol. 4727, pp. 427–442. Springer, Heidelberg (2007)
22. Mangard, S., Oswald, E., Popp, T.: Power Analysis Attacks. Springer, Heidelberg (2007)
23. Massey, J.L.: Guessing and Entropy. In: The proceedings of the IEEE International Symposium on Information Theory, Trondheim, Norway, p. 204 (June 1994)
24. Micali, S., Reyzin, L.: Physically observable cryptography. In: Naor, M. (ed.) TCC 2004. LNCS, vol. 2951, pp. 278–296. Springer, Heidelberg (2004)
25. Petit, C., Standaert, F.-X., Pereira, O., Malkin, T.G., Yung, M.: A Block Cipher based PRNG Secure Against Side-Channel Key Recovery. In: ASIACCS 2008, Tokyo, Japan, pp. 56–65 (March 2008)
26. Pietrzak, K.: A Leakage-Resilient Mode of Operation. In: The proceedings of Eurocrypt 2009, Cologne, Germany. LNCS (April 2009) (to appear)
27. Schindler, W., Lemke, K., Paar, C.: A stochastic model for differential side channel cryptanalysis. In: Rao, J.R., Sunar, B. (eds.) CHES 2005. LNCS, vol. 3659, pp. 30–46. Springer, Heidelberg (2005)
28. Shannon, C.E.: A Mathematical Theory of Communication. Bell System Technical Journal 27, 379–423, 623–656 (1948)
29. Shannon, C.E.: Communication theory of secrecy systems. Bell System Technical Journal 28, 656–715 (1949)
30. Standaert, F.-X., Malkin, T.G., Yung, M.: A Unified Framework for the Analysis of Side-Channel Key Recovery Attacks (extended version), Cryptology ePrint Archive, Report 2006/139
31. Standaert, F.-X., Peeters, E., Archambeau, C., Quisquater, J.-J.: Towards security limits in side-channel attacks. In: Goubin, L., Matsui, M. (eds.) CHES 2006. LNCS, vol. 4249, pp. 30–45. Springer, Heidelberg (2006)
32. Standaert, F.-X., Archambeau, C.: Using subspace-based template attacks to compare and combine power and electromagnetic information leakages. In: Oswald, E., Rohatgi, P. (eds.) CHES 2008. LNCS, vol. 5154, pp. 411–425. Springer, Heidelberg (2008)
33. Standaert, F.-X., Gierlichs, B., Verbauwhede, I.: Partition vs. Comparison Side-Channel Distinguishers: An Empirical Evaluation of Statistical Tests for Univariate Side-Channel Attacks. In: Lee, P.J., Cheon, J.H. (eds.) ICISC 2008. LNCS, vol. 5461, pp. 253–267. Springer, Heidelberg (2009)