

Towards Easy Leakage Certification

François Durvaux, François-Xavier Standaert.

ICTEAM/ELEN/Crypto Group, Université catholique de Louvain, Belgium.

Abstract. Side-channel attacks generally rely on the availability of good leakage models to extract sensitive information from cryptographic implementations. The recently introduced leakage certification tests aim to guarantee that this condition is fulfilled based on sound statistical arguments. They are important ingredients in the evaluation of leaking devices since they allow a good separation between engineering challenges (how to produce clean measurements) and cryptographic ones (how to exploit these measurements). In this paper, we propose an alternative leakage certification test that is significantly simpler to implement than the previous proposal from Eurocrypt 2014. This gain admittedly comes at the cost of a couple of heuristic (yet reasonable) assumptions on the leakage distribution. To confirm its relevance, we first show that it allows confirming previous results of leakage certification. We then put forward that it leads to additional and useful intuitions regarding the information losses caused by incorrect assumptions in leakage modeling.

1 Introduction

Side-channel attacks are an important threat against the security of modern embedded devices. As a result, the search for efficient approaches to secure cryptographic implementations against such attacks has been an ongoing process over the last 15 years. Sound tools for quantifying physical leakages are a central ingredient for this purpose, since they are necessary to balance the implementation cost of concrete countermeasures with the security improvements they provide. Hence, while early countermeasures came with proposals of security evaluations that were sometimes specialized to the countermeasure, more recent works have investigated the possibility to consider evaluation methods that generally apply to any countermeasure. The unified evaluation framework proposed at Eurocrypt 2009 is a typical attempt in this direction [17]. It suggests to analyze cryptographic implementations with a combination of information theoretic and security metrics. The first ones aim at measuring the (worst-case) information leakage independent of the adversary exploiting it, and are typically instantiated with the Mutual Information (MI). The second ones aim at quantifying how efficiently an adversary can take advantage of this leakage in order to turn it into (e.g.) a key recovery, and are typically instantiated with a success rate.

In this context, an important observation is that most side-channel attacks, and in particular any standard Differential Power Analysis (DPA) attack, require a leakage model [10]. This model usually corresponds to an estimation of the leakage Probability Density Function (PDF), possibly simplified to certain statistical moments. Since the exact distribution of (e.g.) power consumption

or electromagnetic radiation measurements is generally unknown, it raises the problem that any physical security evaluation is possibly biased by model errors. In other words, security evaluations ideally require a perfect leakage model (so that all the information is extracted from the measurements). But in practice models are never perfect, so that the quality of the evaluation may highly depend on the quality of the evaluator. This intuition can be captured with the notion of Perceived Information (PI) - that is nothing else than an estimation of the MI biased by the side-channel evaluator's model [15]. Namely, the MI captures the worst-case security level of an implementation, as it corresponds to an (hypothetical) adversary who can perfectly profile the leakage PDF. By contrast, the PI captures its practical counterpart, where actual (statistical) estimation procedures are used by an evaluator, in order to profile the leakage PDF.

Picking up on this problem, Durvaux et al. introduced first “leakage certification” methods at Eurocrypt 2014 [6]. Intuitively, leakage certification starts from the fact that actual leakage models are obtained via PDF estimation, which may lead to both estimation and assumption errors. As a result, and since it seems hard to enforce that such estimated models are perfect, the best that one can hope is to guarantee that they are “good enough”. For estimation errors, this is easily verified using standard cross-validation techniques (in general, estimation errors can anyway be made arbitrarily small by measuring more). For assumption errors, things are more difficult since it requires to find out whether the estimated model is close to an (unknown) perfect model. Interestingly, the Eurocrypt 2014 paper showed that indirect approaches allow determining if this condition is respected, essentially by comparing the model errors caused by incorrect assumptions to estimation errors. That is, let us assume that an evaluator is given a set of leakage measurements to quantify the security of a leaking implementation. As long as the assumption errors measured from these traces remain small in front of the estimation errors, the evaluator is sure that any improvement of his (possibly imperfect) assumptions will not lead to noticeable degradations of the estimated security level – since the impact of these improved assumptions will essentially be hidden by the estimation errors. By contrast, once the assumption errors become significant in front of estimation ones, it means that an improved model is required to extract all the information from the measurements. Hence, leakage certification allows ensuring that the modeling part of an evaluation is sound (i.e. only depends on the implementation – not the evaluator).

In practice, the leakage certification test in [6] requires a number of technical ingredients. Namely, the evaluator first has to characterize the leakages of the target implementation with a sampled (cumulative) distance distribution, and to characterize his model with a simulated (cumulative) distance distribution. Working with distances allows exploiting a univariate goodness-of-fit test even for leakages of large dimensionalities (i.e. it allows comparing the univariate distances between multivariate leakages rather than comparing the multivariate leakages directly). The Cramér-von-Mises divergence is used as a comparison tool in the Eurocrypt 2014 paper. Qualitatively, large divergences between the sampled and simulated distributions essentially mean that the assumptions are

imperfect. Quantitatively, the evaluator then has to determine whether such divergences are significant, by verifying whether they can be explained by assumption errors. This essentially requires computing the p-values when testing the hypothesis that the estimated model is correct (which again requires computing many simulated cumulative distance distributions). Summarizing, the beauty of this approach lies in the fact that it only relies on non-parametric estimations and requires no assumptions on the underlying leakage distributions. But this also comes at the cost of quite computationally intensive tools.

In this paper, we analyze solutions to mitigate the latter drawback, by investigating whether (computationally) cheaper and (conceptually) simpler certification procedures can be obtained at the cost of mild assumptions on the statistical distributions in hand. Two natural options directly come to mind for this purpose, that both aim to avoid dealing with the (expensive to characterize) cumulative leakage distributions directly. One possibility is to “summarize” the leakage distribution with its MI/PI estimates (since they can be used as good indicators of the side-channel security level, as now proven in [5]). Another one is to analyze this distribution “moments by moments”, motivated by the recent results in [13]. In both cases, and following the approach in [6], the main idea remains to compare actual leakage samples generated by a leaking implementation with hypothetical ones generated with the evaluator’s model. Surprisingly, we show that the first approach cannot work, essentially because of situations where model errors in one statistical moment (e.g. the mean) are reflected in another statistical moment (e.g. the variance), which typically arises when using the popular stochastic models in [16], and actually corresponds to the context of epistemic noise discussed in [9]. More interestingly, we also show that a moment-based approach provides excellent results under reasonable assumptions, and can borrow from the “leakage detection tests” that are already used by evaluation laboratories [8, 11]. The resulting leakage certification method is significantly faster than the Eurocrypt 2014 one (and allows reproducing its experiments). We also show that it easily generalizes to masked implementations, and enables extracting very useful intuitions on the origin of the leakages. Eventually, our new tools additionally lead to simple heuristics to approximate the information loss due to incorrect leakage models, which remained an open problem in [6]. Summarizing, we simplify leakage certification into a set of easy-to-implement procedures, hopefully more attractive for evaluation laboratories.

Cautionary note. This paper is about *leakage certification*, which is a different problem than the *leakage detection* one discussed in [8, 11] (despite we indeed borrow some tools from leakage detection to simplify leakage certification). In this respect, Goodwill et al.’s non specific t-test is a natural approach to leakage detection, and allows determining if there is “some” leakage in an implementation, independent of whether it can be exploited (e.g. how many traces do you need to attack). By contrast, leakage certification aims to guarantee that a leakage model that can be exploited in an attack (and, e.g. can be used to determine a key recovery success rate) is close enough to the true leakage model. That is, it aims to make evaluators confident that their attacks are close enough to

the worst-case ones. So leakage detection and certification are essentially complementary. Note that leakage models (and certification) are needed in any attempt to connect side-channel analysis with cryptographic security guarantees (e.g. in leakage resilience [1, 7]), where we will always need an accurate evaluation of the security level, or to build security graphs such as introduced in [20].

2 Background

2.1 Measurement setup

Our experiments are based on measurements of an AES Furious implementation¹ run by an 8-bit Atmel AVR (ATMega644P) microcontroller at a 20 MHz clock frequency. We monitored the voltage variations across a 22 Ω resistor introduced in the supply circuit of our target chip. Acquisitions were performed using a Lecroy WaveRunner HRO 66 oscilloscope running at 625 MHz and providing 8-bit samples. In practice, our evaluations focused on the leakage of the first AES master key byte (but would apply identically to any other enumerable target). Leakage traces were produced according to the following procedure. Let x and s be our target input plaintext byte and subkey, and $y = x \oplus s$. For each of the 256 values of y , we generated 1000 encryption traces, where the rest of the plaintext and key was random (i.e. we generated 256 000 traces in total, with plaintexts of the shape $p = x||r_1||\dots||r_{15}$, keys of the shape $\kappa = s||r_{16}||\dots||r_{30}$, and the r_i 's denoting uniformly random bytes). In order to reduce the memory cost of our evaluations, we only stored the leakage corresponding to the 2 first AES rounds (as the dependencies in our target byte $y = x \oplus s$ typically vanish after the first round, because of the strong diffusion properties of the AES). In the following, we will denote the 1000 encryption traces obtained from a plaintext p including the target byte x under a key κ including the subkey s as: $\text{AES}_{\kappa_s}(p_x) \rightsquigarrow l_y^i$ (with $i \in [1; 1000]$). Eventually, whenever accessing the points of these traces, we will use the notation $l_y^i(\tau)$ (with $\tau \in [1; 10\ 000]$, typically). These subscripts and superscripts will be omitted when not necessary or clear from the context.

2.2 PDF estimation methods

Side-channel attacks such as the standard DPA described in [10] require a leakage model. In general, such models correspond to estimations of the leakage PDF (possibly simplified to certain statistical moments). In the following, we will consider two important PDF estimation techniques for this purpose.

Gaussian templates. The Template Attack (TA) in [3] approximates the leakages using a set of normal distributions. It assumes that each intermediate computation generates Gaussian-distributed samples. In our typical scenario where the targets follow a key addition, we consequently use: $\hat{\text{Pr}}_{\text{model}}[l_y|s, x] \approx \hat{\text{Pr}}_{\text{model}}[l_y|s \oplus x] \sim \mathcal{N}(\mu_y, \sigma_y^2)$, where the ‘‘hat’’ notation is used to denote the

¹ Available at <http://point-at-infinity.org/avraes/>.

estimation of a statistic. This approach requires estimating the sample means and variances for each value $y = x \oplus s$ (and mean vectors / covariance matrices in case of multivariate attacks). We denote the construction of such a model with $\hat{\text{Pr}}_{\text{model}}^{\text{ta}} \leftarrow \mathcal{L}_Y^p$, where \mathcal{L}_Y^p is a set of N_p traces used for profiling.

Regression-based models. To reduce the data complexity of the profiling, an alternative approach proposed by Schindler et al. is to exploit Linear Regression (LR) [16]. In this case, a stochastic model $\hat{\theta}(y)$ is used to approximate the leakage function and built from a linear basis $\mathbf{g}(y) = \{\mathbf{g}_0(y), \dots, \mathbf{g}_{B-1}(y)\}$ chosen by the adversary/evaluator (usually $\mathbf{g}_i(y)$ are monomials in the bits of y). Evaluating $\theta(\hat{y})$ boils down to estimating the coefficients α_i such that the vector $\hat{\theta}(y) = \sum_i \alpha_i \mathbf{g}_i(y)$ is a least-square approximation of the measured leakages L_y . In general, an interesting feature of such models is that they allow trading profiling efforts for online attack complexity, by adapting the basis $\mathbf{g}(y)$. That is, a simpler model with fewer parameters will converge for smaller values of N_p , but a more complex model can potentially approximate the real leakage function more accurately. Compared to Gaussian templates, another feature of this approach is that only a single variance (or covariance matrix) is estimated for capturing the noise (i.e. it relies on an assumption of homoscedastic errors). Again, we denote the constructions of such a model with $\hat{\text{Pr}}_{\text{model}}^{\text{lr}} \leftarrow \mathcal{L}_Y^p$.

2.3 Evaluation metrics

In this subsection, we recall a couple of useful evaluation metrics that have been introduced in previous works on side-channel attacks and countermeasures.

Correlation coefficient. In view of the popularity of the Correlation Power Analysis (CPA) distinguisher in the literature [2], a natural candidate evaluation metric is Pearson’s correlation coefficient. In the non-profiled setting, an a-priori (e.g. Hamming weight) model is used for computing the metric. The evaluator then estimates the correlation between his measured leakages and the modeled leakages of a target intermediate value. In our AES example, it would lead to $\hat{\rho}(L_Y(\tau), \text{model}_{\text{cpa}}(Y))$. In practice, this estimation is performed by sampling (i.e. measuring) N_t test traces from the leakage distribution (we denote the set of these N_t test traces as \mathcal{L}_Y^t). Next, and in order to avoid possible biases due to an incorrect a-priori choice of leakage model, a natural solution is to extend the previous proposal to the profiled setting. In this case, the evaluator will start by estimating a model from N_p profiling traces: $\hat{\text{model}}_{\text{cpa}} \leftarrow \mathcal{L}_Y^p$ (with $\mathcal{L}_Y^p \perp \mathcal{L}_Y^t$). In practice, $\hat{\text{model}}_{\text{cpa}}$ can be seen as a simplification of the previous Gaussian templates, that only includes estimates for the first-order moments of the leakages. That is, for any time sample τ , we have $\hat{\text{model}}_{\text{cpa}}(y) = \hat{m}_y^1(\tau) = \hat{\text{E}}_i(L_y^i(\tau))$, with \hat{m}_y^1 a first-order moment and $\hat{\text{E}}$ the sample mean operator.

Mutual and perceived information. In theory, the worst-case security level of an implementation can be measured with a MI metric. Taking advantage of the notations in Section 2.1 and considering the standard case where a key byte S is targeted, it amounts to estimate the following quantity:

$$\text{MI}(S; X, L) = \text{H}[S] + \sum_{s \in \mathcal{S}} \text{Pr}[s] \sum_{x \in \mathcal{X}} \text{Pr}[x] \sum_{l_y^i \in \mathcal{L}_t} \text{Pr}_{\text{chip}}[l_y^i | s, x] \cdot \log_2 \text{Pr}_{\text{chip}}[s | x, l_y^i].$$

When summing over all s and x values, and a sufficiently large number of leakages, the estimation tends to the correct MI. Yet, as mentioned in introduction, the chip distribution $\text{Pr}_{\text{chip}}[l_y^i | s, x]$ is generally unknown to the evaluator. So in practice, the best that we can hope is to compute the following PI:

$$\hat{\text{PI}}(S; X, L) = \text{H}[S] + \sum_{s \in \mathcal{S}} \text{Pr}[s] \sum_{x \in \mathcal{X}} \text{Pr}[x] \sum_{l_y^i \in \mathcal{L}_t} \text{Pr}_{\text{chip}}[l_y^i | s, x] \cdot \log_2 \hat{\text{Pr}}_{\text{model}}[s | x, l_y^i],$$

where $\hat{\text{Pr}}_{\text{model}} \leftarrow \mathcal{L}_Y^p$ is typically obtained using the previous Gaussian templates or LR-based models. Under the assumption that the model is properly estimated, it is shown in [10] that the CPA and PI metrics are essentially equivalent in the context of standard univariate side-channel attacks (i.e. exploiting a single leakage point $l_y^i(\tau)$ at a time). By contrast, only the PI naturally extends to multivariate attacks. It can be interpreted as the amount of information leakage that will be exploited by an adversary using an estimated model. So just as the MI is a good predictor for the success rate of an ideal TA exploiting the perfect model Pr_{chip} , the PI is a good predictor for the success rate of an actual TA exploiting the “best available” model $\hat{\text{Pr}}_{\text{model}}$ obtained thanks to profiling.

Moments-correlating DPA. Eventually, and in order to extend the CPA distinguisher to higher-order moments, the Moments-Correlating Profiled DPA (MCP-DPA) has been introduced in [13]. It features essentially the same steps as a profiled CPA. The only difference is that the adversary first estimates d th-order statistical moments with his profiling traces, and then uses $\hat{\text{model}}_{\text{mcp-dpa}}^d(y) = \hat{m}_y^d(\tau)$, with \hat{m}_y^d a d th-order moment. For concreteness, we will consider d ’s up to four (i.e. the sample mean for $d = 1$, variance for $d = 2$, skewness for $d = 3$ and kurtosis for $d = 4$), which allows us discussing the relevant case-study of a masked implementation with two shares. Yet, the tool naturally extends to any d . One useful feature of this distinguisher is that it embeds the same “metric” intuition as CPA: the higher the correlation estimated with MCP-DPA, the more efficient the corresponding attack exploiting a moment of given order.

2.4 Estimating a metric with cross-validation

Estimating a metric α from a leaking implementation holds in two steps. First, a model has to be estimated from a set of profiling traces \mathcal{L}_p : $\text{model} \leftarrow \mathcal{L}_p$. Second, a set of test traces \mathcal{L}_t (following the true distribution Pr_{chip}) is used to estimate the metric: $\hat{\alpha} \leftarrow (\mathcal{L}_t, \text{model})$. As a result, two main types of errors can arise. First, the number of traces in the profiling set may be too low to estimate the model accurately (which corresponds to estimation errors). Second, the model may not be able to accurately predict the distribution of samples in the test set, even after intensive profiling (which then corresponds to assumption errors).

In order to verify that estimations in a security evaluation are sufficiently accurate, the solution used in [6] is to exploit cross-validation. In general, this technique allows gauging how well a predictive (here leakage) model performs in practice. For k -fold cross-validations, the set of evaluation traces \mathcal{L} is first split into k (non overlapping) sets $\mathcal{L}^{(i)}$ of approximately the same size. Let us define the profiling sets $\mathcal{L}_p^{(j)} = \bigcup_{i \neq j} \mathcal{L}^{(i)}$ and the test sets $\mathcal{L}_t^{(j)} = \mathcal{L} \setminus \mathcal{L}_p^{(j)}$. The sample metric is then repeatedly computed k times for $1 \leq j \leq k$ as follows. First, we build a model from a profiling set: $\text{model}^{(j)} \leftarrow \mathcal{L}_p^{(j)}$. Then we estimate the metric with the associated test set $\hat{\alpha}^{(j)} \leftarrow (\mathcal{L}_t^{(j)}, \text{model}^{(j)})$. Cross-validation protects evaluators from obtaining too optimistic sample metric values due to over-fitting, since the test computations are always performed with an independent data set. Finally, the k outputs can be averaged in order to get an unbiased metric estimate, and their spread characterizes the result’s accuracy.

3 A motivating negative result

As mentioned in introduction, detecting assumption errors is generally more challenging than detecting estimation errors (which is easily done with the previous cross-validation). Intuitively, it requires to investigate the likelihood that samples obtained from a leaking device can indeed be explained by an estimated model, which requires a (multivariate) goodness-of-fit test. Since such tests are computationally intensive, an appealing alternative would be to check whether the samples obtained from the leaking device lead to a PI that is at least close enough to the MI: this would guarantee a good estimation of the security level. But we again face the problem that the MI is unknown, which imposes trying indirect approaches. That is, we would need an metric counterpart to the sampled simulated distance distribution in [6], which would typically correspond to the following Hypothetical (mutual) Information (HI):

$$\hat{\text{HI}}(S; X, L) = \text{H}[S] - \sum_{s \in \mathcal{S}} \text{Pr}[s] \sum_{x \in \mathcal{X}} \text{Pr}[x] \sum_{l_y^i \in \mathcal{L}_t} \hat{\text{Pr}}_{\text{model}}[l_y^i | s, x] \cdot \log_2 \hat{\text{Pr}}_{\text{model}}[s | x, l_y^i].$$

Intuitively, this HI corresponds to the amount of information that would be extracted from an hypothetical implementation that would exactly leak according to the model $\hat{\text{Pr}}_{\text{model}}$. In itself, the HI is useless to the evaluator, as it is actually disconnected from the chip distribution. For example, even a totally incorrect model (i.e. leading to a negative PI) would lead to a positive HI. By contrast, we could hope that as long as the HI and PI are “close”, the assumption errors are “small enough” for the number of measurements considered in the security evaluation. Furthermore, we could use a simple hypothesis test to detect non-closeness. For a number of traces N in the evaluation set, this would require to compute estimates $\hat{\text{PI}}(S; X, L)^{(j)}$ and $\hat{\text{HI}}(S; X, L)^{(j)}$ with cross-validation, and to check whether these estimates come from different (univariate) distributions. If they significantly differ, we would conclude that the model exhibits assumption errors that degrade the estimated security level, in a similar fashion as in [6].

Unfortunately, and despite it can detect certain assumption errors, this approach cannot succeed in general. A simple counter-example can be explained in the context of LR. Say an adversary estimates a model with a linear basis, which leads to significant differences between the actual (mean) leakages and the ones suggested by the model. Then, because of the homoscedastic error assumption, the single variance of the LR-based model will reflect this error (i.e. capture both physical noise and model error). As a result, whenever this type of error increases, the PI will decrease (as expected) but the HI will also decrease (contrary to the MI). So testing the consistency between the PI and HI estimates will not reveal the inconsistencies between the PI estimates and the true MI.

4 A new method to detect assumption errors

Despite negative, the previous counter-example suggests two interesting tracks for simplifying leakage certification tests. First, summarizing a complete distribution into representative metrics (e.g. such as the PI) allows taking advantage of simpler statistical tests. Second, since the fact that the homoscedastic errors assumption is not fulfilled implies that errors made in the estimation of certain statistical moments (or more generally, parameters) of a distribution are reflected in other statistical moments of this distribution, a natural approach is to test the relevance of a model “moment by moment”. That is, for a number of traces N in an evaluation set, one could verify that the moments estimated from actual leakage samples are hard to tell apart from the moments estimated from the model (with the same number of samples N). Based on this idea, our simplified method to detect assumption errors will be based on the following two hypotheses (one strictly necessary and the other optional but simplifying).

1. *The leakage distribution is well represented by its statistical moments.* This corresponds to the classical “moment problem” in statistics, for which there exist counter-examples (e.g. the log-normal distribution is not uniquely characterized by its moments). So our (informal) assumption is that these counter-examples will not be significant for our experimental case-studies.
2. *The sampled estimates of our statistical moments are approximately Gaussian-distributed.* This directly derives from the central limit theorem and actually depends on the number of samples used in the estimations (which will become sufficient as the leakages become more noisy, e.g. in the case of protected implementations that are most relevant for concrete investigations).

Let us add a couple of words of motivation for those assumptions. First recall that we know from the previous results in [6] that leakage certification is possible without such assumptions, at the cost of somewhat involved statistical reasoning and estimations. So it seems natural to investigate alternative (heuristic) paths allowing to reach similar conclusions. As will be shown next, this is indeed the case of our simplified approach for a couple of relevant scenarios. Second, statistical moments are at the core of the reasoning regarding the masking countermeasure. That is, the security order of an implementation is generally defined

as the lowest informative moment in the leakage distribution (minus one) – see [5] for an extensive discussion of this issue. Besides, many concrete (profiled and non-profiled) side-channel attacks are based (implicitly or explicitly) on parametric PDF estimation techniques that rely on the estimation of moments (e.g. the Gaussian templates and LR-based models in Section 2.2, but also second-order attacks such as [4, 14]). So an approach based on an analysis of moments seems well founded in these cases.² As a result, and maybe most importantly, we believe that the following tools open interesting research avenues regarding the intuitive evaluation of leaking devices based on their moments.

As for the Gaussian assumption, our motivation is even more pragmatic, and relates to the observation that simple t-tests are becoming de facto standards in the preliminary evaluation of leaking devices [8, 11]. So we find it appealing to rely on statistical tools that are already widespread in the CHES community, and to connect them with leakage certification. As will be clear next, this allows us to use the same evaluation method for statistical moments of different orders. However, we insist that it is perfectly feasible to refine our approach by using a well adapted test for each statistical moment (e.g. F-test for variances, ...).

4.1 Test specification

The main idea behind our new leakage certification method is to compare (actual) d th-order moments \hat{m}_y^d estimated from the leakages with (simulated) d th-order moments \tilde{m}_y^d estimated from the evaluator’s model $\hat{\text{Pr}}_{\text{model}}$ (by sampling this model). Thanks to our second assumption, this comparison can simply be performed based on Student’s t-test. For this purpose, we need multiple estimations of the moments \hat{m}_y^d and \tilde{m}_y^d , that we will obtain thanks to an approach inspired from Section 2.4 (although there is no cross-validation involved here).

More precisely, we start by splitting the full set of evaluation traces \mathcal{L} into k (non overlapping) sets of approximately the same size $\mathcal{L}^{(j)}$, with $1 \leq j \leq k$. From these k subsets, we produce k estimates of (actual) d^{th} -order moments $\hat{m}_y^{d,(j)}$, each of them from a set $\mathcal{L}^{(j)}$. We then produce a set of simulated traces $\tilde{\mathcal{L}}$ that has the same size and corresponds to the same intermediate values as the real evaluation set \mathcal{L} , but where the leakages are sampled according to the model that we want to evaluate. In other words, we first build the model $\hat{\text{Pr}}_{\text{model}} \leftarrow \mathcal{L}$, and then generate a simulated set of traces $\tilde{\mathcal{L}} \leftarrow \hat{\text{Pr}}_{\text{model}}$. Based on $\tilde{\mathcal{L}}$, we produce k estimates of (simulated) d^{th} -order moments $\tilde{m}_y^{d,(j)}$, each of them from a set $\tilde{\mathcal{L}}^{(j)}$, as done for the real set of evaluation traces. From these real and simulated moments estimates, we compute the following quantities:

$$\begin{aligned} \hat{\mu}_y^d &= \hat{\text{E}}_j(\hat{m}_y^{d,(j)}), & \hat{\sigma}_y^d &= \sqrt{\hat{\text{var}}_j(\hat{m}_y^{d,(j)})}, \\ \tilde{\mu}_y^d &= \hat{\text{E}}_j(\tilde{m}_y^{d,(j)}), & \tilde{\sigma}_y^d &= \sqrt{\hat{\text{var}}_j(\tilde{m}_y^{d,(j)})}, \end{aligned}$$

² Non-parametric PDF estimations do not suffer from assumption errors (at the cost of a significantly increased estimation cost), so are out of scope here.

where $\hat{\text{var}}$ is the sample variance operator. Eventually, we simply estimate the t statistic (next denoted with Δ_y^d) as follows:

$$\Delta_y^d = \frac{\hat{\mu}_y^d - \tilde{\mu}_y^d}{\sqrt{\frac{(\hat{\sigma}_y^d)^2 + (\tilde{\sigma}_y^d)^2}{k}}}.$$

The p -value of this t statistic within the associated Student’s distribution returns the probability that the observed difference is the result of estimations issues, and is computed as:

$$p = 2 \times (1 - \text{CDF}_t(|\Delta_y^d|, d_f)),$$

where CDF_t is the Student’s t cumulative distribution function, and d_f is its number freedom degrees.³ In other words, a small p -value indicates that the model is incorrect with high probability. Concretely, the only parameter to set in this test is the number of non overlapping sets k . Following [6], we used $k = 10$ which is a rather standard value in the literature. Note that increasing k has very limited impact on the accuracy of our conclusions since all variance estimates in the t -test are normalized by k . By contrast it increases the time complexity of the test (so keeping k reasonably small is in general a good strategy).

5 Simulated experiments

In order to validate our moment-based certification method, we first analyze a couple of simulated experiments, where we can control the assumption errors. In particular, and in order to keep these simulations reasonably close to concrete attacks, we consider four distinct scenarios. In the first one (reported in Figure 1), both the leakage function Pr_{chip} and the leakage model $\hat{\text{Pr}}_{\text{model}}$ follow a Gaussian distribution, but the model’s estimated mean differs from the true distribution. In the second one (reported in Figure 2), the leakage function and the leakage model again follow a Gaussian distribution, this time with a model error on the variance. These two examples informally correspond to the context of LR-based attacks, where the basis is not large enough to capture the exact mean values. Our third and fourth examples correspond to a slightly different setting aimed to emulate a masked implementation, for which the true distribution is typically a mixture [19]. So in the third scenario (reported in Figure 3), the leakage function has a Gaussian mixture distribution with a non-zero skewness, while the leakage model is still a Gaussian approximation. And in the fourth scenario (reported in Figure 4), the leakage function has a Gaussian mixture distribution with a non-zero kurtosis, while the leakage model is still a Gaussian approximation. In all cases, we represent the true distribution and the biased model, the estimated moments for these two distributions, and the p -value of our certification test.

³ Student’s t distribution is a parametric probability density function whose only parameter is its number of freedom degrees, that can be directly derived from k and the previous σ estimates as: $d_f = (k - 1) \times [(\hat{\sigma}_y^d)^2 + (\tilde{\sigma}_y^d)^2] / [(\hat{\sigma}_y^d)^4 + (\tilde{\sigma}_y^d)^4]$.

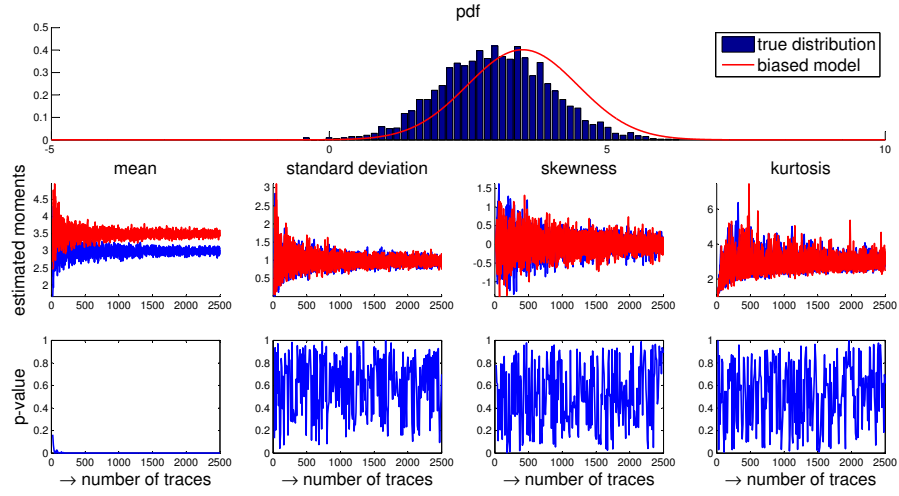


Fig. 1. Gaussian leakages, Gaussian model, error in the estimated mean.

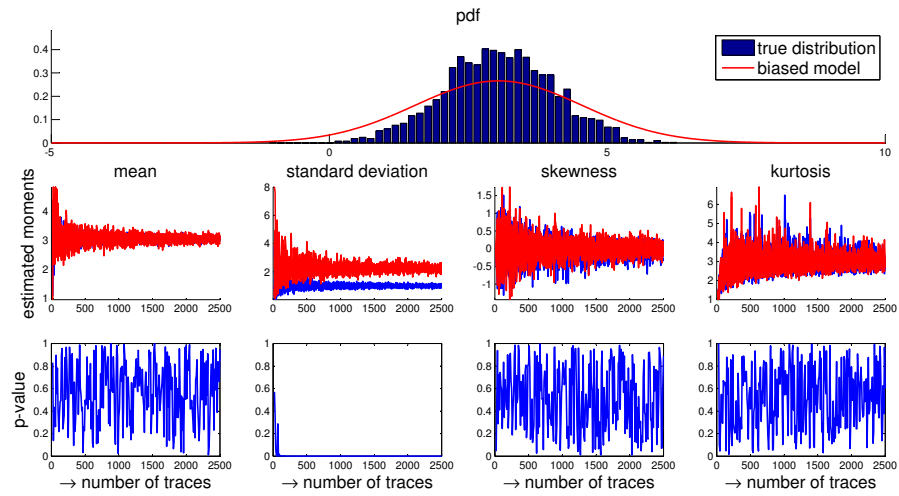


Fig. 2. Gaussian leakages, Gaussian model, error in the estimated variance.

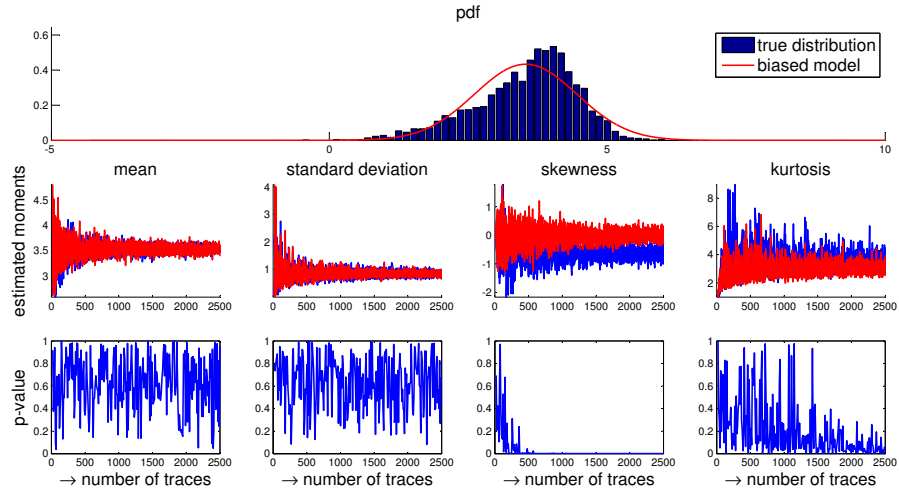


Fig. 3. Gaussian mixture leakages, Gaussian model, error in the estimated skewness.

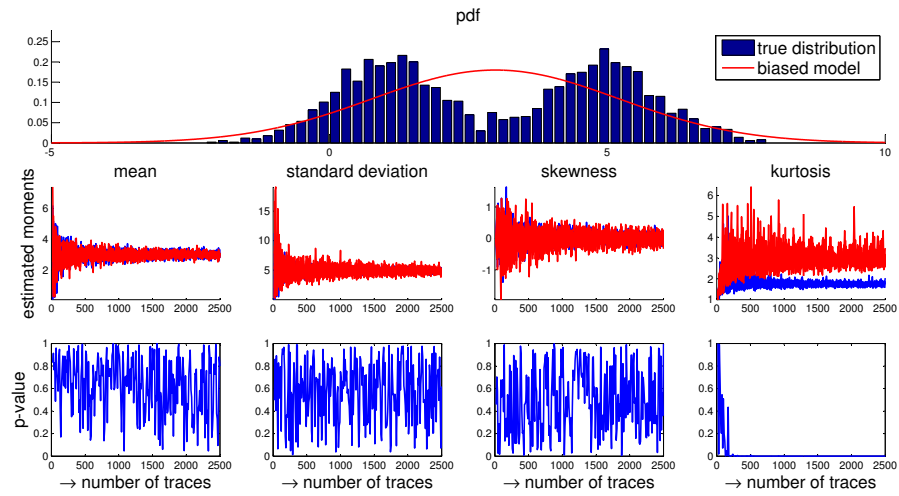


Fig. 4. Gaussian mixture leakages, Gaussian model, error in the estimated kurtosis.

The results are mostly as expected and confirm the simplicity of the method. That is, as the number of measurements in the evaluation set increases, we are able to detect the assumption errors in all cases. The only difference between the applications to different moments is that errors on higher-order moments may be more difficult to detect as the noise increases. This difference is caused by the same argument that justifies the relevance of the higher-order masking countermeasure. Namely, the sampling complexity when estimating the moments of a distribution increases exponentially in d . However, this is not a limitation of the certification test: if such errors are not detected for a given evaluation set, it just means that their impact is still small in front of assumption errors at this stage of the evaluation. Besides, we note that the respective relevance of the model errors on different moments will be further discussed in Section 7.

6 Measured experiments

In order to obtain a fair comparison with the results provided in [6], we also applied our new leakage certification method to the same case-study. That is, we used the measurement setup from Section 2.1 and evaluated the relevance of two important profiling methods, namely the Gaussian TA and LR, for the most informative time sample in our leakage traces (i.e. with maximum PI).

The main difference with the previous simulated experiments is that we now have to test 256 models independently (each of them corresponding to a target intermediate value $y = x \oplus s$). Our results are represented in Figure 5, where we plot the p-values output by our different t-tests in greyscale, for four statistical moments (i.e. the mean, variance, skewness and kurtosis). A look at the first two moments essentially confirms the results of Durvaux et al. More precisely, the Gaussian templates seem to capture the measured leakages quite accurately (for the 256,000 traces in our evaluation set). By contrast, the linear regression quickly exhibits inconsistencies. Interestingly, assumption errors appear both in the means and in the variances, which corresponds to the expected intuition. That is, errors in the means are detected because for most target intermediate values, the actual leakage cannot be accurately predicted by a linear combination of the S-box output bits. And errors in the variances appear because the LR-based models rely on the homoscedastic error assumption and capture both physical noise and noise due to assumption errors in a single term.

By contrast, and quite intriguingly, a look at the last two moments (i.e. skewness and kurtosis) also shows some differences with the results in [6]. That is, we remark that even for Gaussian templates, small model errors appear in these higher-order moments. This essentially corresponds to the fact that our measured leakages do not have perfectly key-independent skewness and kurtosis, as we assume in Gaussian PDF estimations. This last observation naturally raises the question whether these errors are significant, i.e. do they contradict the results of the Eurocrypt 2014 leakage certification test? In the next section, we show that it is not the case, and re-conciliate both approaches by investigating the respective informativeness of the four moments in our new test.

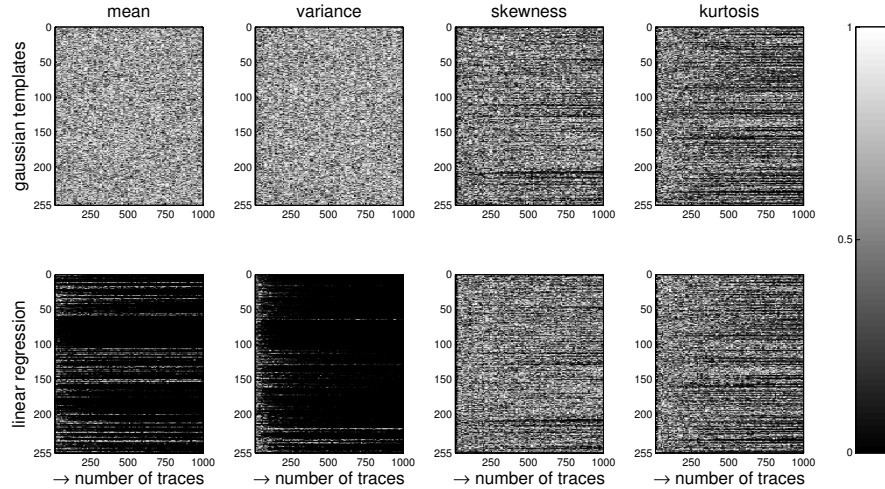


Fig. 5. Results of the new leakage certification test for actual measurements.

7 Quantifying the information loss

Since Figure 5 suggests the existence of (small) model errors in our Gaussian templates, that are due to an incorrect characterization of the third- and fourth-order moments in our leakage traces, we now want to investigate whether these errors are leading to significant information losses. Fortunately, our “per-moment” approach to leakage certification also allows simple investigations in this direction (which incidentally and heuristically answers one of the open questions in [6], about the information loss due to model errors). In particular, we can simply use the MCP-DPA mentioned in Section 2.3 for this purpose. Roughly, this tool computes the correlation between a simplified model (that corresponds to d th-order moments of the leakage distribution) to samples raised to the power d (possibly centered or standardized if we consider centered and standardized moments). As discussed in [13], the resulting estimated correlation coefficient features a “metric intuition”: the higher the value of the MCP-DPA distinguisher computed for an order d , the more efficient the MCP-DPA attack exploiting this statistical order of the actual leakage distribution. Hence, computing the value of the MCP-DPA distinguisher for different values of d should solve our problem.

Concretely, we start by applying MCP-DPA in the traditional sense and exploit cross-validation for this purpose, this time following exactly Section 2.4. That is, the set of evaluation traces \mathcal{L} is again split into k (non overlapping) sets $\mathcal{L}^{(i)}$ of approximately the same size, and we use profiling sets $\mathcal{L}_p^{(j)} = \bigcup_{i \neq j} \mathcal{L}^{(i)}$ and test sets $\mathcal{L}_t^{(j)} = \mathcal{L} \setminus \mathcal{L}_p^{(j)}$. We then repeatedly compute the d th-order moments $\hat{m}_y^{d,(j)} \leftarrow \mathcal{L}_p^{(j)}$, and the d th-order MCP-DPA distinguisher:

$$\text{MCP-DPA}^{(j)}(d) = \hat{\rho} \left(\hat{M}_Y^{d,(j)}, (L_y)^d \leftarrow \mathcal{L}_t^{(j)} \right).$$

As previously mentioned, it corresponds to the sample correlation between the random variable representing the estimated moments \hat{M}_Y^d , and the random variable corresponding to the leakage samples coming from the test set $L_y \leftarrow \mathcal{L}_t^{(j)}$, raised to power d (possibly centered or standardized if we consider centered and standardized moments). The $k = 10$ estimates for this MCP-DPA metric are represented in the top part of Figure 6. We additionally considered two slightly tweaked versions of MCP-DPA, where we rather estimate Gaussian TA (resp. LR-based) models $\hat{\text{Pr}}_{\text{model}}^{\text{ta}}$ (resp. $\hat{\text{Pr}}_{\text{model}}^{\text{lr}}$), and consider the two (resp. one) key-dependent moments from these models to compute the metric. These tweaked MCP-DPAs are represented in the middle (resp. lower) part of the figure.

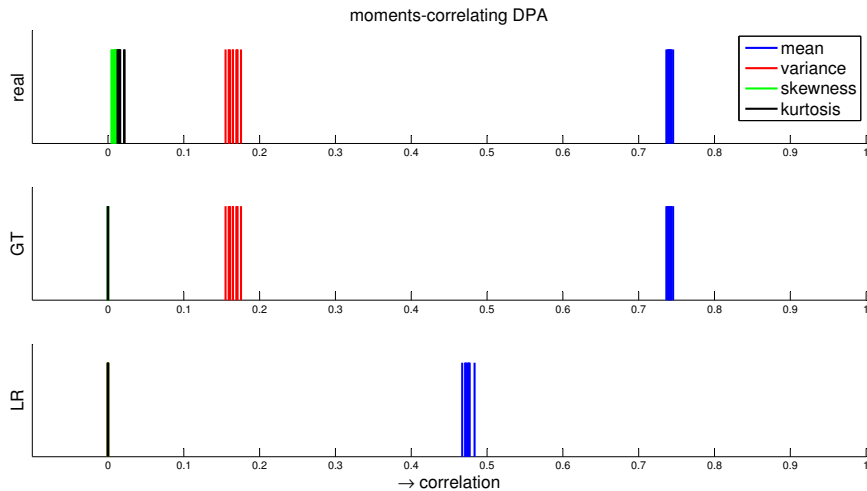


Fig. 6. MCP-DPA results actual measurements.

Our main observations are as follows. First, the upper part of the figure suggests that the most informative moments in our leakage traces are the mean and variance. There is indeed a small amount of information in the skewness and kurtosis. But by considering the classical rule-of-thumb that the measurement complexity of a correlation-based attack is inversely proportional to the square of its correlation coefficient, we can see that the additional information gain in these higher-order moments is very limited in our context. This observation backs up the conclusions of the generic leakage certification test in [6] that Gaussian templates are sufficiently accurate for our evaluation set. It can be similarly quantified by using the following approximation from [10]:

$$I(X, Y) \approx -\frac{1}{2} \times \log_2 \left(1 - \rho(X, Y)^2 \right),$$

considering that attacks exploiting two moments would take advantage of the sum of their respective information (i.e. considering them as independent information channels), and using the formula in [?] which shows that the measurement complexity of a side-channel attack is inversely proportional to the MI/PI leaked by a target implementation. Next, we also see that TA-based and LR-based MCP-DPA yield no information in the higher-order moments, which trivially derives from the fact that they rely on a Gaussian assumption. Eventually, we notice that the information loss between LR-based models and TA-based models can be approximated thanks to the correlation between their moments. For example, and considering the means in Figure 6, we can compute the value of the LR-based MCP-DPA distinguisher – worth ≈ 0.48 in the figure – by multiplying the value of the TA-based MCP-DPA distinguisher – worth ≈ 0.74 in the figure – by $\hat{\rho}(\hat{M}_Y^{d,ta}, \hat{M}_Y^{d,lr})$ – worth ≈ 0.65 in our experiments (i.e. by taking advantage of the “product rule” for the correlation coefficient in [18]).

Those last tools are admittedly informal. Yet, we believe they provide a useful variety of heuristics allowing evaluators to analyze the results of their certification tests. In particular, they lead to easy-to-exploit intuitions regarding the impact (or lack thereof) of model errors detected in moments of a given order. As discussed in the beginning of Section 4, further formalizing these findings, and possibly putting forward relevant scenarios where our simplified approach leads to significant shortcomings, is an interesting scope for further research.

Eventually, we mention that from the time complexity point-of-view, the leakage certification tools in this paper are considerably more efficient than the previous ones from [6]. Strict comparisons are hard to obtain since our current implementations are prototype ones, and further optimizations could certainly be investigated. But roughly speaking, generating leakage certification plots for 256 leakage models as in Figure 5 is completed in minutes of computations on a standard desktop computer, whereas it typically took hours with the Eurocrypt 2014 tools. Since the cost of our heuristic leakage certification method is essentially similar to the one of a CPA, it can easily be applied on full leakage traces, in particular if some high performance computing can be exploited to take advantage of the parallel nature of the certification problem [12].

Acknowledgements. F.-X. Standaert is a research associate of the Belgian Fund for Scientific Research (FNRS-F.R.S.). This work has been funded in parts by the European Commission through the ERC project 280141 (CRASH).

References

1. Sonia Belaïd, Vincent Grosso, and François-Xavier Standaert. Masking and leakage-resilient primitives: One, the other(s) or both? *Cryptography and Communications*, 7(1):163–184, 2015.
2. Eric Brier, Christophe Clavier, and Francis Olivier. Correlation power analysis with a leakage model. In Marc Joye and Jean-Jacques Quisquater, editors, *Cryptographic*

- Hardware and Embedded Systems - CHES 2004: 6th International Workshop Cambridge, MA, USA, August 11-13, 2004. Proceedings*, volume 3156 of *Lecture Notes in Computer Science*, pages 16–29. Springer, 2004.
3. Suresh Chari, Josyula R. Rao, and Pankaj Rohatgi. Template attacks. In Burton S. Kaliski Jr., Çetin Kaya Koç, and Christof Paar, editors, *Cryptographic Hardware and Embedded Systems - CHES 2002, 4th International Workshop, Redwood Shores, CA, USA, August 13-15, 2002, Revised Papers*, volume 2523 of *Lecture Notes in Computer Science*, pages 13–28. Springer, 2002.
 4. Guillaume Dabosville, Julien Doget, and Emmanuel Prouff. A new second-order side channel attack based on linear regression. *IEEE Trans. Computers*, 62(8):1629–1640, 2013.
 5. Alexandre Duc, Sebastian Faust, and François-Xavier Standaert. Making masking security proofs concrete - or how to evaluate the security of any leaking device. In Elisabeth Oswald and Marc Fischlin, editors, *Advances in Cryptology - EUROCRYPT 2015 - 34th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Sofia, Bulgaria, April 26-30, 2015, Proceedings, Part I*, volume 9056 of *Lecture Notes in Computer Science*, pages 401–429. Springer, 2015.
 6. François Durvaux, François-Xavier Standaert, and Nicolas Veyrat-Charvillon. How to certify the leakage of a chip? In Phong Q. Nguyen and Elisabeth Oswald, editors, *Advances in Cryptology - EUROCRYPT 2014 - 33rd Annual International Conference on the Theory and Applications of Cryptographic Techniques, Copenhagen, Denmark, May 11-15, 2014. Proceedings*, volume 8441 of *Lecture Notes in Computer Science*, pages 459–476. Springer, 2014.
 7. Stefan Dziembowski and Krzysztof Pietrzak. Leakage-resilient cryptography. In *49th Annual IEEE Symposium on Foundations of Computer Science, FOCS 2008, October 25-28, 2008, Philadelphia, PA, USA*, pages 293–302. IEEE Computer Society, 2008.
 8. Gilbert Goodwill, Benjamin Jun, Josh Jaffe, and Pankaj Rohatgi. A testing methodology for side channel resistance validation. NIST non-invasive attack testing workshop, 2011. http://csrc.nist.gov/news_events/non-invasive-attack-testing-workshop/papers/08_Goodwill.pdf.
 9. Annelie Heuser, Olivier Rioul, and Sylvain Guilley. Good is not good enough - deriving optimal distinguishers from communication theory. In Lejla Batina and Matthew Robshaw, editors, *Cryptographic Hardware and Embedded Systems - CHES 2014 - 16th International Workshop, Busan, South Korea, September 23-26, 2014. Proceedings*, volume 8731 of *Lecture Notes in Computer Science*, pages 55–74. Springer, 2014.
 10. Stefan Mangard, Elisabeth Oswald, and François-Xavier Standaert. One for all - all for one: unifying standard differential power analysis attacks. *IET Information Security*, 5(2):100–110, 2011.
 11. Luke Mather, Elisabeth Oswald, Joe Bandenburg, and Marcin Wójcik. Does my device leak information? an a priori statistical power analysis of leakage detection tests. In Kazue Sako and Palash Sarkar, editors, *Advances in Cryptology - ASIACRYPT 2013 - 19th International Conference on the Theory and Application of Cryptology and Information Security, Bengaluru, India, December 1-5, 2013, Proceedings, Part I*, volume 8269 of *Lecture Notes in Computer Science*, pages 486–505. Springer, 2013.
 12. Luke Mather, Elisabeth Oswald, and Carolyn Whitnall. Multi-target DPA attacks: Pushing DPA beyond the limits of a desktop computer. In Palash Sarkar and Tetsu

- Iwata, editors, *Advances in Cryptology - ASIACRYPT 2014 - 20th International Conference on the Theory and Application of Cryptology and Information Security, Kaoshiung, Taiwan, R.O.C., December 7-11, 2014. Proceedings, Part I*, volume 8873 of *Lecture Notes in Computer Science*, pages 243–261. Springer, 2014.
13. Amir Moradi and François-Xavier Standaert. Moments-correlating DPA. *IACR Cryptology ePrint Archive*, 2014:409, 2014.
 14. Emmanuel Prouff, Matthieu Rivain, and Régis Bevan. Statistical analysis of second order differential power analysis. *IEEE Trans. Computers*, 58(6):799–811, 2009.
 15. Mathieu Renaud, François-Xavier Standaert, Nicolas Veyrat-Charvillon, Dina Kamel, and Denis Flandre. A formal study of power variability issues and side-channel attacks for nanoscale devices. In Kenneth G. Paterson, editor, *Advances in Cryptology - EUROCRYPT 2011 - 30th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Tallinn, Estonia, May 15-19, 2011. Proceedings*, volume 6632 of *Lecture Notes in Computer Science*, pages 109–128. Springer, 2011.
 16. Werner Schindler, Kerstin Lemke, and Christof Paar. A stochastic model for differential side channel cryptanalysis. In Josyula R. Rao and Berk Sunar, editors, *Cryptographic Hardware and Embedded Systems - CHES 2005, 7th International Workshop, Edinburgh, UK, August 29 - September 1, 2005, Proceedings*, volume 3659 of *Lecture Notes in Computer Science*, pages 30–46. Springer, 2005.
 17. François-Xavier Standaert, Tal Malkin, and Moti Yung. A unified framework for the analysis of side-channel key recovery attacks. In Antoine Joux, editor, *Advances in Cryptology - EUROCRYPT 2009, 28th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Cologne, Germany, April 26-30, 2009. Proceedings*, volume 5479 of *Lecture Notes in Computer Science*, pages 443–461. Springer, 2009.
 18. François-Xavier Standaert, Eric Peeters, Gaël Rouvroy, and Jean-Jacques Quisquater. An overview of power analysis attacks against field programmable gate arrays. *Proceedings of the IEEE*, 94(2):383–394, 2006.
 19. François-Xavier Standaert, Nicolas Veyrat-Charvillon, Elisabeth Oswald, Benedikt Gierlichs, Marcel Medwed, Markus Kasper, and Stefan Mangard. The world is not enough: Another look on second-order DPA. In Masayuki Abe, editor, *Advances in Cryptology - ASIACRYPT 2010 - 16th International Conference on the Theory and Application of Cryptology and Information Security, Singapore, December 5-9, 2010. Proceedings*, volume 6477 of *Lecture Notes in Computer Science*, pages 112–129. Springer, 2010.
 20. Nicolas Veyrat-Charvillon, Benoît Gérard, and François-Xavier Standaert. Security evaluations beyond computing power. In Thomas Johansson and Phong Q. Nguyen, editors, *Advances in Cryptology - EUROCRYPT 2013, 32nd Annual International Conference on the Theory and Applications of Cryptographic Techniques, Athens, Greece, May 26-30, 2013. Proceedings*, volume 7881 of *Lecture Notes in Computer Science*, pages 126–141. Springer, 2013.