



HAL
open science

Optimal side-channel attacks for multivariate leakages and multiple models

Nicolas Bruneau, Sylvain Guilley, Annelie Heuser, Marion Damien, Olivier
Rioul

► **To cite this version:**

Nicolas Bruneau, Sylvain Guilley, Annelie Heuser, Marion Damien, Olivier Rioul. Optimal side-channel attacks for multivariate leakages and multiple models. PROOFS 2016 Security Proofs for Embedded Systems, Aug 2016, Santa Barbara, United States. 10.1007/s13389-017-0170-9 . hal-02300060

HAL Id: hal-02300060

<https://telecom-paris.hal.science/hal-02300060>

Submitted on 12 Aug 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Optimal Side-Channel Attacks for Multivariate Leakages and Multiple Models

Nicolas Bruneau^{1,2}, Sylvain Guilley^{1,3}, Annelie Heuser^{1*},
Damien Marion^{1,3}, and Olivier Rioul^{1,4}

¹ LTCI, CNRS, Télécom ParisTech, Université Paris-Saclay, 75 013 Paris, France

² STMicroelectronics, AST Division, 13 790 Rousset, France

³ Secure-IC S.A.S., Threat Analysis Business Line, 35 510 Cesson-Sévigné, France

⁴ CMAP, École polytechnique, Université Paris-Saclay, 91 128 Palaiseau, France

firstname.lastname@telecom-paristech.fr

Abstract. Side-channel attacks allow to extract secret keys from embedded systems like smartcards or smartphones. In practice, the side-channel signal is measured as a trace consisting of several samples. Also, several sensitive bits are manipulated in parallel, each leaking differently. Therefore, the informed attacker needs to devise side-channel distinguishers that can handle both *multivariate* leakages and *multiple* models. In the state-of-the-art, these two issues have two independent solutions: on the one hand, dimensionality reduction can cope with multivariate leakage; on the other hand, online stochastic approach can cope with multiple models.

In this paper, we combine both solutions to derive closed-form expressions of the resulting *optimal* distinguisher in terms of matrix operations, in all situations where the model can be either profiled offline or regressed online. Optimality here means that the success rate is maximized for a given number of traces. We recover known results for uni- and bi-variate models (including correlation power analysis), and investigate novel distinguishers for multiple models with more than two parameters. In addition, following ideas from the AsiaCrypt'2013 paper "*Behind the Scene of Side-Channel Attacks*", we provide fast computation algorithms in which the traces are accumulated prior to computing the distinguisher values.

1 Introduction

Side-channel attacks allow to extract secret keys from cryptographic devices. Template attacks [3] have been introduced as the strongest analysis method. They consist in two phases: (i) a profiling offline phase where the leakage model of the device under attack is characterized; (ii) an attack online phase in which the secret key is extracted using fresh measurements along with the pre-characterized

* Annelie Heuser is a Google European fellow in the field of privacy and is partially founded by this fellowship.

model. Such attacks are known to use a maximum likelihood principle to ensure the highest possible success probability (see, eg., [8]).

In this paper we study *optimal* attacks with the best possible success probability when extracting the secret key¹. We leverage on such optimal distinguishers to answer the following question: how to attack with the best probability of success when the *leakage* is *multivariate* and the *model* are *multiple*? An initial empirical work has already been carried out in [13] and confirmed that this type of approach can be very fruitful².

Contributions. We derive closed-form expressions for the optimal distinguishers in all situations where the model is known (e.g., using profiling) or regressed online. In the case of a known univariate model, we recover the results in [2]. However, our “fully matrix” formalism makes equations simpler and proofs shorter. Moreover, compared to [2] we extend the leakage model to the case where the traces are not necessarily centered, thereby allowing a more natural application on real traces. In the realistic “stochastic attack” situation where the model is *parametric*, i.e. where the coefficients of the model are unknown, we express the optimal distinguisher by maximizing success over the whole set of possible coefficients. Finally, we provide fast computation algorithms for our novel distinguishers, which happen to be remarkably simple and efficient.

Outline. The remainder of this paper is organized as follows. Sec. 2 provides a modelization of a side-channel attack that is generic enough to capture many different multivariate scenarios. The main results of this paper are outlined in Sec. 3. Sec. 4 presents experimental results on simulated traces and real-world acquisition campaigns. Conclusions and perspectives are in Sec. 5.

2 Notations and Leakage Model

2.1 Notations

We let X denote the leakage measurements, Y the model, N the measurement noise, and α the link between the model and the measurements³. The model Y

¹ The success probability in key recovery is chosen as a figure of merit for optimization. Such an objective is typical of “pure” side-channel attacks. Other approaches [7, 10, 15] relax the condition that the key found by the side-channel analysis be ranked first and complements it with a key enumeration stage. This yields a data vs. complexity tradeoff that is not explored in this paper.

² *Multi-target attacks* [10, 16] have a somewhat different goal, namely the best aggregation of information about several subparts of a key, possibly leaking at different times with different models, in order to recover the full key efficiently. Here we consider only *one* multivariate leakage model and focus on recovering *one* subpart of the key. However, our derivation is capable of handling multivariate leakages and models and may still be combined with the multi-target approaches.

³ Notations X, Y are consistent with the usual convention in machine learning, where X is for the collected data and Y for the classification labels.

depends on a key guess k , an n -bit vector, and on some known text T (usually also an n -bit vector) e.g., through a function ϕ such that $Y = \phi(T, k)$. A well-known example is $Y = w_H(T \oplus k)$, where w_H is the Hamming weight function. However, in general, some parameters of the model are unknown. To remain generic, we do not detail further the link between Y and (T, k) . As it is customary in side-channel analysis, the correct key is denoted by k^* . The corresponding model using the correct key $Y(k^*)$ is denoted by Y^* .

Let Q be the number of queries (number of measurements), D be the data dimensionality (number of time samples per measurement trace) and S be the model dimensionality ($\phi : \mathbb{F}_2^n \times \mathbb{F}_2^n \rightarrow \mathbb{R}^S$ is a vectorial function, with S components). Roman letters in **bold** indicate vectors or matrices that have a dimension in Q , i.e., which are different for each trace $q = 1, 2, \dots, Q$. More precisely, \mathbf{X} represents the full attack campaign, a matrix of $D \times Q$ measurement samples. The q -th trace is denoted X_q which is a $D \times 1$ column vector. Similarly, for the q -th trace, the $S \times 1$ column vector Y_q represents the deterministic part of the model while the $D \times 1$ column vector N_q is the corresponding measurement noise with $D \times D$ correlation matrix Σ .

We denote by $\text{tr}(\cdot)$ the trace of a square matrix, that is the sum of its diagonal terms. Note that $\text{tr}(AB) = \text{tr}(BA)$ for compatible matrix dimensions. Let $\|\cdot\|_2$ denote the Euclidean norm of a $1 \times Q$ row vector. Thus $\|\mathbf{X}\|_2^2 = \mathbf{X}\mathbf{X}^\top = \text{tr}(\mathbf{X}^\top\mathbf{X})$, where $(\cdot)^\top$ is the transposition operator. Finally let $\|\cdot\|_F$ denote the Frobenius norm of a matrix (square root of the sum of its squared elements), such that $\|M\|_F^2 = \text{tr}(MM^\top)$.

2.2 General Model

We make the following simplifying assumptions. First, the (environmental) noise is steady, e.g., chip temperature and supply voltage do not vary during the side-channel measurements. Thus N_1, N_2, \dots, N_Q are independent and identically distributed (i.i.d.) (denoted by N with index q dropped). Second, the attacker does not inject partial information gathered from the leakage analysis into a possible choice of plaintexts/ciphertexts (nonadaptive attack)⁴. Thus Y_1, Y_2, \dots, Y_Q are assumed i.i.d. (denoted by Y). Under the adopted leakage model it follows that the leakage measurements X_1, X_2, \dots, X_Q are also i.i.d. (denoted by X).

A distinguisher \mathcal{D} maps a collection of leakages \mathbf{x} and texts \mathbf{t} to an estimation of the secret key k^* . Let us recall that \mathbf{x} and \mathbf{t} are realizations of random variables \mathbf{X} and \mathbf{T} : \mathbf{x} is a $D \times Q$ matrix of real numbers (the acquisition campaign) and \mathbf{t} is a $1 \times Q$ vector of n -bit words (bytes when $n = 8$) which are the publicly known plaintext or ciphertext bytes. An *optimal* distinguisher maximizes the probability of success $\mathcal{D}(\mathbf{x}, \mathbf{t}) = k^*$.

The simplest situation occurs when X consists in a modulation of Y plus noise, in which case we let α be the signal envelope. In real traces, however, we face the more general situation where the model can be offset by some quantity

⁴ In fact, our results tolerate chosen texts, but consider them as observed inputs to the attack. We do not optimize the attack according to chosen inputs.

the general case being an S -dimensional *parametric* model with $S \geq 2$ components. For this reason, we consider α as a $D \times S$ matrix and we set in matrix notation

$$\mathbf{X} = \alpha \mathbf{Y}^* + \mathbf{N} \quad (1)$$

where \mathbf{X} is $D \times Q$, α is $D \times S$, \mathbf{Y}^* is $S \times Q$, and \mathbf{N} is $D \times Q$. Notice that our convention to consider traces as lines and dimensions as rows allows us to write the deterministic part of the leakage as $\alpha \mathbf{Y}^*$ which writes more naturally than the opposite order where traces would be viewed as a vertical time series.

For each trace $q = 1, 2, \dots, Q$, we assume that the vector $N = N_q$ follows that same multivariate normal distribution $\mathcal{N}(0, \Sigma)$, where the $D \times D$ correlation matrix $\Sigma = \mathbb{E}(NN^T)$ is assumed known to the attacker⁵. Since Σ is assumed symmetric positive definite, there exists a matrix $\Sigma^{1/2}$, which is such that $\Sigma^{1/2} \Sigma^{1/2} = \Sigma$. We refer to $\Sigma^{1/2}$ as the *standard deviation noise* matrix.

The model (1) used throughout the paper is quite generic and have multiple facets depending on the choice of S and the respective values given to α and \mathbf{Y} . This is discussed next.

2.3 Examples with $S = 2$ and $S = 9$

For $S = 1$, the traces consist only in a *modulation* of the model plus noise as in [2]. When considering traces that are not only modulated but also have an *offset* term we have $S = 2$. We then write the 2-dimensional model as $\begin{pmatrix} \mathbf{Y} \\ \mathbf{1} \end{pmatrix}$, where \mathbf{Y} and $\mathbf{1}$ are $1 \times Q$ matrices (Y_1, Y_2, \dots, Y_Q) and $(1, 1, \dots, 1)$. The $D \times 2$ matrix α in (1) actually takes the special form $\begin{pmatrix} \alpha & \beta \end{pmatrix}$ where β is the offset value.

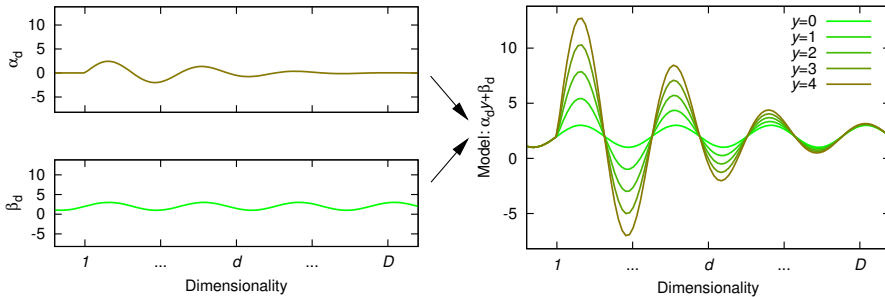


Fig. 1. Example of leakage model with $S = 2$ and a model in Hamming weight, with $n = 4$ values (no noise is added)

An illustration is provided in Fig. 1 where the parameter $\beta \in \mathbb{R}^D$ is the waveform when there is no signal, whereas $\alpha \in \mathbb{R}^D$ is the signal envelope. The

⁵ We may simplify (2) by incorporating $\beta \mathbf{1}$ into the noise expectation, but the noise is intrinsically zero-mean and it is clearer to exhibit a specific offset term.

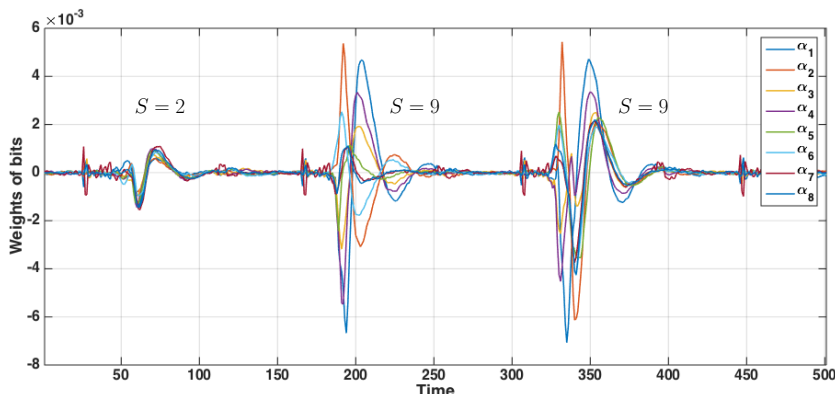


Fig. 2. Leakage evaluation of traces from DPA contest V4 (knowing the mask)

complete model is the sum $\alpha Y + \beta$, where Y is the Hamming weight of some intermediate variable (such as the XOR operation $T \oplus k$) on $n = 4$ bits. While the leakage signal may be represented as a continuous curve as illustrated in Fig. 1, the practical acquisition is done through a temporal series of D “discrete samples”, typically within one clock period. For $S = 2$, we thus write (1) as

$$\mathbf{X} = \alpha \mathbf{Y}^* + \beta \mathbf{1} + \mathbf{N} \quad (2)$$

where \mathbf{X} is $D \times Q$, α and β are $D \times 1$, \mathbf{Y}^* and $\mathbf{1} = (1, \dots, 1)$ are $1 \times Q$, and \mathbf{N} is $D \times Q$. Here \mathbf{Y} is assumed centered: $\mathbb{E}(\mathbf{Y}) = \mathbf{0} = (0, \dots, 0)$ (since the non-centered part is captured by the $\beta \mathbf{1}$ term) and of unit variance for every q : $\text{Var}(Y_q) = \mathbb{E}(Y_q^2) = 1$.

For $S \geq 2$, the actual value of S reflects the complexity of the model. For example, in the *weighted sum of bits* model, the model for each trace can be written as $\sum_{s=1}^n \alpha_s Y_s + \beta$, where Y_s is the s th bit of the n -bit sensitive variable Y . Accordingly, we have

$$S = n + 1, \quad \alpha = (\alpha_1 \dots \alpha_n \beta), \quad \mathbf{Y} = (\mathbf{Y}_1 \dots \mathbf{Y}_n \mathbf{1})^\top. \quad (3)$$

This leakage model is more complex than before but may arise in practice. For example, Fig. 2 plots the coefficients $\alpha_1, \dots, \alpha_8$ estimated of the traces taken from an ATMega smartcard—the datasets are available from the DPA contest V4 team [14]. In particular one can observe that samples around $[50, 80]$ are ordered by Hamming weight: this part of the trace resembles the upper left part of Fig. 1 for $S = 2$. By analysing the $(n+1)$ -variate model of (3), one can indeed see that around $[50, 80]$, the vectors $\alpha_1, \dots, \alpha_8$ are almost identical. However, samples in intervals $[170, 250]$ or $[330, 400]$ have a more complex model. These times, the eight vectors $\alpha_1, \dots, \alpha_8$ are clearly different, the leakage is 9-variate.

In the sequel, we consider both types of attacks: those with *offline profiling* where α for each component of the model is precharacterized like in Fig. 2

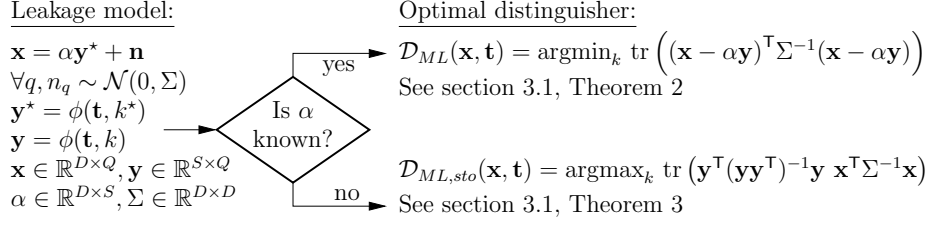


Fig. 3. Mathematical expression for multivariate ($D \geq 1$) optimal attacks with a linear combination of models ($S \geq 1$)

and also those where the model is *learned online* like in a Linear Regression Attack [5].

3 Theoretical Results and Implementation

3.1 General Mathematical Expressions

In this section we derive the mathematical expression of the optimal distinguisher \mathcal{D} in the general case of multivariate leakage ($D \geq 1$), and multiple models ($S \geq 1$). An illustration of our results is given in Fig. 3 for the case when the leakage is completely known (or profiled as in the template attack) and when the leakage is unknown and estimated online.

Definition 1 (Optimal Distinguisher Knowing or Ignoring α).

$$\mathcal{D}_{ML}(\mathbf{x}, \mathbf{t}) = \operatorname{argmax}_{k \in \mathbb{F}_2} p(\mathbf{x}|\mathbf{t}) \quad \text{and} \quad \mathcal{D}_{ML,sto}(\mathbf{x}, \mathbf{t}) = \operatorname{argmax}_{k \in \mathbb{F}_2} \max_{\alpha \in \mathbb{R}^{D \times S}} p(\mathbf{x}|\mathbf{t}, \alpha).$$

In both cases (Theorems 2 and 3 below) the result is a distinguisher which is computed using simple matrix operations. While \mathcal{D}_{ML} resembles a template attack with Gaussian templates [3], $\mathcal{D}_{ML,sto}$ is a novel expression that results from a non-trivial maximization over the matrix α and may be interpreted as a generalization of a multivariate correlation power attack [1].

Theorem 2. *The optimal maximum likelihood (ML) distinguisher [8] for Gaussian noise writes*

$$\mathcal{D}_{ML}(\mathbf{x}, \mathbf{t}) = \operatorname{argmin}_k \operatorname{tr} \left((\mathbf{x} - \alpha \mathbf{y})^\top \Sigma^{-1} (\mathbf{x} - \alpha \mathbf{y}) \right). \quad (4)$$

Proof. From [8] we have $\mathcal{D}_{ML}(\mathbf{x}, \mathbf{t}) = \operatorname{argmax}_k p(\mathbf{x}|\mathbf{y})$ while from (1) we see that $p(\mathbf{x}|\mathbf{y}) = p_{\mathbf{N}}(\mathbf{x} - \alpha \mathbf{y})$. From the i.i.d. assumption the noise density $p_{\mathbf{N}}(\mathbf{n})$

is given by

$$p_{\mathbf{N}}(\mathbf{n}) = \prod_{q=1}^Q \frac{1}{\sqrt{(2\pi)^D |\det \Sigma|}} \exp -\frac{1}{2} n_q^\top \Sigma^{-1} n_q \quad (5)$$

$$= \frac{1}{(2\pi)^{DQ/2}} \frac{1}{(\det \Sigma)^{Q/2}} \exp -\frac{1}{2} \left(\sum_{q=1}^Q n_q^\top \Sigma^{-1} n_q \right) \quad (6)$$

$$= \frac{1}{(2\pi)^{DQ/2} (\det \Sigma)^{Q/2}} \exp -\frac{1}{2} \text{tr} (\mathbf{n}^\top \Sigma^{-1} \mathbf{n}). \quad (7)$$

Thus $p_{\mathbf{N}}(\mathbf{x} - \alpha \mathbf{y})$ is maximum when the expression $\text{tr} (\mathbf{n}^\top \Sigma^{-1} \mathbf{n})$ for $\mathbf{n} = \mathbf{x} - \alpha \mathbf{y}$ is minimum. \square

Theorem 3. *The optimal stochastic multivariate attack is given by*

$$\mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t}) = \underset{k \in \mathbb{F}_2^n}{\text{argmax}} \text{tr} (\mathbf{y}^\top (\mathbf{y}\mathbf{y}^\top)^{-1} \mathbf{y} \mathbf{x}^\top \Sigma^{-1} \mathbf{x}) \quad (8)$$

for which the optimal value of α is given by

$$\alpha^{\text{opt}} = (\mathbf{x}\mathbf{y}^\top)(\mathbf{y}\mathbf{y}^\top)^{-1}. \quad (9)$$

For the proof, we need some known results of linear algebra (Lemma 4) and linear regression (Lemma 5).

Lemma 4. *Let \mathbf{b} an $S \times Q$ matrix, with $S < Q$. The $S \times S$ matrix $\mathbf{b}\mathbf{b}^\top$ is invertible if and only if \mathbf{b} has full rank S , i.e., if and only if the S lines of \mathbf{b} are independent.*

Proof. For any $S \times 1$ column vector x , $x^\top \mathbf{b}\mathbf{b}^\top x = \|\mathbf{b}^\top x\|^2 = 0$ implies $\mathbf{b}^\top x = 0$ hence $x = 0$. Hence the matrix $\mathbf{b}\mathbf{b}^\top$ is positive definite. \square

Lemma 5. *Let \mathbf{a} , \mathbf{b} and α be respectively $1 \times Q$, $S \times Q$ and $1 \times S$ with $S < Q$, where \mathbf{b} has full rank S . Then $\|\mathbf{a} - \alpha \mathbf{b}\|_2$ reaches its minimum for $\alpha = \mathbf{a}\mathbf{b}^\top (\mathbf{b}\mathbf{b}^\top)^{-1}$.*

Proof. Expanding the squared norm gives $\|\mathbf{a} - \alpha \mathbf{b}\|_2^2 = (\mathbf{a} - \alpha \mathbf{b})(\mathbf{a} - \alpha \mathbf{b})^\top = \mathbf{a}\mathbf{a}^\top - 2\alpha \mathbf{b}\mathbf{a}^\top + \alpha \mathbf{b}\mathbf{b}^\top \alpha^\top$. Therefore, the gradient $\frac{\partial}{\partial \alpha} \|\mathbf{a} - \alpha \mathbf{b}\|_2^2 = -2\mathbf{b}\mathbf{a}^\top + 2\mathbf{b}\mathbf{b}^\top \alpha^\top$ vanishes if and only if $\alpha^\top = (\mathbf{b}\mathbf{b}^\top)^{-1} \mathbf{b}\mathbf{a}^\top$, i.e., $\alpha = \mathbf{a}\mathbf{b}^\top (\mathbf{b}\mathbf{b}^\top)^{-1}$ where we have used the fact that $\mathbf{b}\mathbf{b}^\top$ is invertible by Lemma 4. \square

Proof (Proof of Theorem 3). Let $\mathbf{x}' = \Sigma^{-1/2} \mathbf{x}$ and $\mathbf{y}' = (\mathbf{y}\mathbf{y}^\top)^{-1/2} \mathbf{y}$. The optimal distinguisher minimizes the following expression over $\alpha \in \mathbb{R}^{D \times S}$:

$$\text{tr} \left((\mathbf{x} - \alpha \mathbf{y})^\top \Sigma^{-1} (\mathbf{x} - \alpha \mathbf{y}) \right) = \text{tr} \left((\mathbf{x}' - \alpha' \mathbf{y}') (\mathbf{x}' - \alpha' \mathbf{y}')^\top \right) = \sum_{d=1}^D \|\mathbf{x}'_d - \alpha'_d \mathbf{y}'_d\|^2.$$

By Lemma 5 the minimization over α'_d yields $\alpha'_d = (\mathbf{x}'_d \mathbf{y}^\top)(\mathbf{y} \mathbf{y}^\top)^{-1}$ for all $d = 1, \dots, D$. This gives $\alpha' = (\mathbf{x}' \mathbf{y}^\top)(\mathbf{y} \mathbf{y}^\top)^{-1}$ hence $\alpha = (\mathbf{x} \mathbf{y}^\top)(\mathbf{y} \mathbf{y}^\top)^{-1}$, which remarkably does *not* depend on Σ .

The minimized value of the distinguisher is thus

$$\begin{aligned} \min_{\alpha} \operatorname{tr} \left((\mathbf{x} - \alpha \mathbf{y})^\top \Sigma^{-1} (\mathbf{x} - \alpha \mathbf{y}) \right) &= \operatorname{tr} \left((\mathbf{x} - \alpha^{\text{opt}} \mathbf{y})^\top \Sigma^{-1} (\mathbf{x} - \alpha^{\text{opt}} \mathbf{y}) \right) \\ &= \operatorname{tr} \left((\text{Id} - \mathbf{y}^\top (\mathbf{y} \mathbf{y}^\top)^{-1})^2 \mathbf{x}^\top \Sigma^{-1} \mathbf{x} \right) \\ &= \operatorname{tr} \left(\mathbf{x}^\top \Sigma^{-1} \mathbf{x} \right) - \operatorname{tr} \left(\mathbf{y}^\top (\mathbf{y} \mathbf{y}^\top)^{-1} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} \right) \end{aligned}$$

where Id is the $D \times D$ identity matrix and where $\operatorname{tr} \left(\mathbf{x}^\top \Sigma^{-1} \mathbf{x} \right)$ is a constant independent of k . This proves Theorem 3. \square

The expression of $\mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t})$ given in Theorem 3 consists in the trace of a $Q \times Q$ matrix, which can be admittedly very large. It can be, however, rewritten in a way that is easier to compute when Q is much greater than S and D :

Corollary 6 (Alternative Expression of $\mathcal{D}_{\text{ML,sto}}$). *Letting $\mathbf{x}' = \Sigma^{-1/2} \mathbf{x}$, and $\mathbf{y}' = (\mathbf{y} \mathbf{y}^\top)^{-1/2} \mathbf{y}$ as in the proof of Theorem 3, we have*

$$\mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t}) = \operatorname{argmax}_{k \in \mathbb{F}_2^n} \|\mathbf{x}' \mathbf{y}'^\top\|_F. \quad (10)$$

Here the Frobenius norm is of a $D \times S$ matrix.

Proof. Let us write $(\mathbf{y} \mathbf{y}^\top)^{-1} = (\mathbf{y} \mathbf{y}^\top)^{-1/2} (\mathbf{y} \mathbf{y}^\top)^{-1/2}$ in (8). By the properties of the trace,

$$\begin{aligned} \operatorname{tr} \left(\mathbf{y}^\top (\mathbf{y} \mathbf{y}^\top)^{-1} \mathbf{y} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} \right) &= \operatorname{tr} \left(\underbrace{(\mathbf{y} \mathbf{y}^\top)^{-1/2} \mathbf{y} (\Sigma^{-1/2} \mathbf{x})^\top}_{S \times D} \underbrace{\left((\mathbf{y} \mathbf{y}^\top)^{-1/2} \mathbf{y} (\Sigma^{-1/2} \mathbf{x}) \right)^\top}_{D \times S} \right) \\ &= \operatorname{tr} \left((\mathbf{y}' \mathbf{x}'^\top) (\mathbf{y}' \mathbf{x}'^\top)^\top \right) = \|\mathbf{x}' \mathbf{y}'^\top\|_F^2. \quad \square \end{aligned}$$

3.2 Mathematical Expressions for $S = 2$

In order to provide interpretations for the optimal distinguisher expressions, we detail how an optimal attack unfolds when the leakage consists in a sum of a modulated scalar model and an offset ($S = 2$). The cases for profiled attacks (denoted $\mathcal{D}_{\text{ML}}^{S=2}$) and non-profiled attacks (denoted $\mathcal{D}_{\text{ML,sto}}^{S=2}$) are presented in Fig. 4.

Interestingly, when $S = 2$, the template attack can decompose in two steps (affine projection followed by a Euclidean distance to the model). Remarkably, the projection vector is the same for all key guesses. This extends similar results [2] where only the linear relationship between leakage and model is explored. As for the online attack, $\mathcal{D}_{\text{ML,sto}}^{S=2}$ consists in a sum of square of CPA attacks on transformed data, aiming at orthogonalizing the noise.

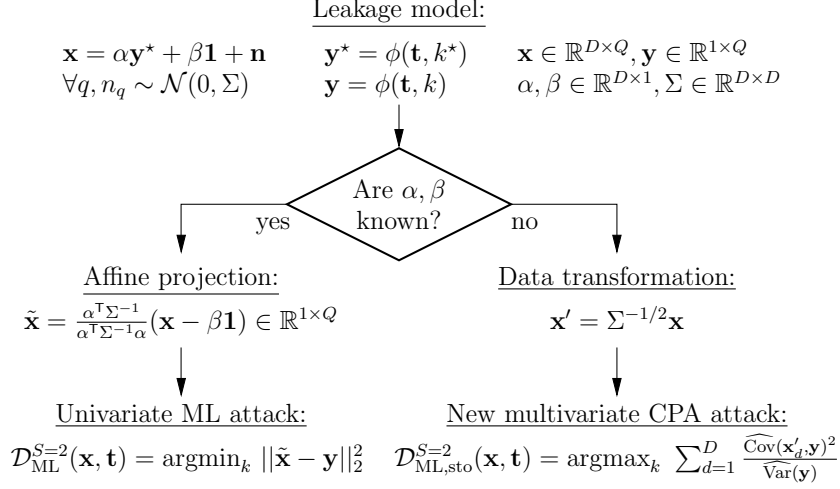


Fig. 4. Modus operandi for multivariate ($D \geq 1$) optimal attacks with one model \mathbf{Y} associated to envelope $\alpha \in \mathbb{R}^{D \times 1}$ and a constant offset $\beta \in \mathbb{R}^{D \times 1}$ ($S = 2$)

3.3 Efficient Implementation

Both \mathcal{D}_{ML} and $\mathcal{D}_{\text{ML,sto}}$ can be optimized using the idea presented in [9]. This article applies a precomputation step in the case the number of traces is larger than the number of possible plaintexts ($Q > \#\mathcal{T} = 2^n$). In this case, all summations \sum_q can be advantageously replaced by $\sum_t \sum_{t_q=t}$. In most cases, the sum $\sum_{t_q=t}$ can be achieved on the fly, and does not involve an hypothesis on the key. Therefore, a speed gain of 2^n (the cardinality of the key space) is expected.

Such optimization strategy can be applied to \mathcal{D}_{ML} . Indeed, let us define $\mathbf{x}' = \Sigma^{-1/2} \mathbf{x}$ and $\alpha' = \Sigma^{-1/2} \alpha$. Then,

$$\begin{aligned} \mathcal{D}_{\text{ML}}(\mathbf{x}, \mathbf{t}) &= \operatorname{argmin}_k \sum_{d=1}^D \|\mathbf{x}'_d - \alpha'_d \mathbf{y}\|_2^2 \quad (\text{see Corollary 6}) \\ &= \operatorname{argmin}_k \sum_{d=1}^D \sum_{t \in \mathbb{F}_2^n} \left(\sum_{q/t_q=t} x'^2_{d,q} - 2 \sum_{q/t_q=t} x'_{d,q} \alpha'_d y(t, k) + \left(\sum_{q/t_q=t} 1 \right) (\alpha'_d y(t, k))^2 \right) \\ &= \operatorname{argmin}_k \sum_{d=1}^D \sum_{t \in \mathbb{F}_2^n} -2 \left(\underbrace{\sum_{q/t_q=t} x'_{d,q}}_{\text{denoted as } x'_{d,t}} \right) \alpha'_d y(t, k) + \left(\underbrace{\sum_{q/t_q=t} 1}_{\text{denoted as } n_t} \right) (\alpha'_d y(t, k))^2. \quad (11) \end{aligned}$$

Notice that at line (11), the term $\sum_{q/t_q=t} x'^2_{d,q}$ which does not depend on the key, is simplified. The fast version of this computation is given in Alg. 1.

```

input :  $\mathbf{x}, \mathbf{t}$ 
output:  $\mathcal{D}_{\text{ML}}(\mathbf{x}, \mathbf{t})$ 

// Initialize to zero a matrix  $x'_{d,t}$  of size  $D \times 2^n$ 
// Initialize to zero a vector  $n_t$  of length  $2^n$ 
1 for  $q \in \{1, \dots, Q\}$  do
2 |  $x'_{t_q} \leftarrow x'_{t_q} + \Sigma^{-1/2} x_q$ 
3 |  $n_{t_q} \leftarrow n_{t_q} + 1$ 
4 end
5 return  $\text{argmin}_{k \in \mathcal{K}} \sum_{d=1}^D \sum_t -2x'_t \alpha'_d y(t, k) + n_t (\alpha'_d y(t, k))^2$  // As in (11)

```

Algorithm 1: Fast computation algorithm for \mathcal{D}_{ML}

The same optimization applies to $\mathcal{D}_{\text{ML,sto}}$. Indeed, in the expression (10) of $\mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t}) = \text{argmax}_k \|\mathbf{x}' \mathbf{y}'^\top\|_F^2$, one can write

$$\|\mathbf{x}' \mathbf{y}'^\top\|_F^2 = \sum_{s,d} \left(\sum_{q=1}^Q x'_{d,q} y'_{s,q} \right)^2 = \sum_{s,d} \left(\underbrace{\sum_{t \in \mathbb{F}_2^n} \left(\sum_{q/t_q=t} x'_{d,t} \right)}_{\text{denoted as } x'_{d,t}} \underbrace{\left(y'_s(t, k) \right)}_{\text{denoted as } y'_{s,t}} \right)^2. \quad (12)$$

This means that x' can be obtained by simple accumulation. The term $y'_s(t, k)$ requires the computation of $\mathbf{y} \mathbf{y}^\top$. In the case $Q \gg 1$, it can be assumed that the texts \mathbf{t} are uniformly distributed. Hence, when $Q \rightarrow +\infty$, by the law of large numbers,

$$\frac{1}{Q} \mathbf{y} \mathbf{y}^\top = \frac{1}{Q} \sum_{q=1}^Q y_q y_q^\top = \sum_{t \in \mathbb{F}_2^n} \frac{\sum_{q/t_q=t} 1}{Q} y(t, k) y(t, k)^\top \rightarrow \frac{1}{2^n} \sum_{t \in \mathbb{F}_2^n} y(t, k) y(t, k)^\top.$$

Therefore, in (12), $y'_s(t)$ can also be precomputed. To the best of our knowledge, this optimization has never been discussed previously. The resulting distinguishing procedure is given in Alg. 2. At line (4), the argument of the Frobenius norm can be computed by a fast matrix multiplication. Also, we notice that the matrix inversion in Alg. 2 is actually a precomputation which involves only the model. Hence the computational complexity of the optimal stochastic attack simply consists in traces accumulation per class, and as many matrix products and Frobenius norms as keys to be guessed.

4 Practical Results

4.1 Characterization of Σ

In this article, we assume that the attacker knows the noise covariance matrix. We give a straightforward procedure for the estimation.

```

input :  $\mathbf{x}, \mathbf{t}$ 
output:  $\mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t})$ 

// Precompute the  $\#\mathcal{K} = 2^n$  matrices  $y'(k)$  of size  $S \times 2^n$ , such that
 $y'(k) = (\frac{1}{2^n} \sum_t y(t, k) y(t, k)^\top)^{-1/2} y(k)$ . Note that there is only one
matrix if the EIS holds [12, Def. 2]

// Initialize to zero a matrix  $x'_{d,t}$  of size  $D \times 2^n$ 
1 for  $q \in \{1, \dots, Q\}$  do
2 |  $x'_{t_q} \leftarrow x'_{t_q} + \Sigma^{-1/2} x_q$  // In-place accumulation of a row in matrix  $x'$ 
3 end
4 return  $\text{argmax}_{k \in \mathcal{K}} \|x' y'(k)^\top\|_F$  // As in (12)
    
```

Algorithm 2: Fast computation algorithm for $\mathcal{D}_{\text{ML,sto}}$ when \mathbf{t} is balanced

1. collect Q traces (i.e., a matrix $\mathbf{x} \in \mathbb{R}^{D \times Q}$) where the plaintext is fixed to a given value,
2. estimate Σ as $\hat{\Sigma} = \frac{1}{Q-1} \left(\mathbf{x} - \frac{1}{Q} \mathbf{x} \mathbf{1}^\top \mathbf{1} \right) \left(\mathbf{x} - \frac{1}{Q} \mathbf{x} \mathbf{1}^\top \mathbf{1} \right)^\top$, where $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^{1 \times Q}$. This estimator is *sample covariance matrix*, which is unbiased.

Remark 1. Notice that Σ cannot be obtained by a direct profiling on the same traces to be used for the attack. Indeed, in those traces, the plaintext is varying, hence the attacker would use for $\hat{\Sigma}$ the covariance matrix of $\mathbf{x} - \alpha^{\text{opt}} \mathbf{y}$, where α^{opt} is equal to $\alpha^{\text{opt}} = (\mathbf{x} \mathbf{y}^\top) (\mathbf{y} \mathbf{y}^\top)^{-1}$ (recall (9)). Hence, $\hat{\Sigma} = \frac{1}{Q-1} (\mathbf{x} - \alpha^{\text{opt}} \mathbf{y}) (\mathbf{x} - \alpha^{\text{opt}} \mathbf{y})^\top$. But the distinguisher $\mathcal{D}_{\text{ML,sto}}$ is

$$\begin{aligned}
 \mathcal{D}_{\text{ML,sto}}(\mathbf{x}, \mathbf{t}) &= \underset{k \in \mathbb{F}_2^n}{\text{argmin}} \min_{\alpha \in \mathbb{R}^{D \times S}} \text{tr} \left((\mathbf{x} - \alpha \mathbf{y})^\top \hat{\Sigma}^{-1} (\mathbf{x} - \alpha \mathbf{y}) \right) \\
 &= \underset{k \in \mathbb{F}_2^n}{\text{argmin}} \min_{\alpha \in \mathbb{R}^{D \times S}} \text{tr} \left(\hat{\Sigma}^{-1} (\mathbf{x} - \alpha \mathbf{y}) (\mathbf{x} - \alpha \mathbf{y})^\top \right) \\
 &= \underset{k \in \mathbb{F}_2^n}{\text{argmin}} \text{tr} \left(\hat{\Sigma}^{-1} (\mathbf{x} - \alpha^{\text{opt}} \mathbf{y}) (\mathbf{x} - \alpha^{\text{opt}} \mathbf{y})^\top \right) \tag{13}
 \end{aligned}$$

$$\begin{aligned}
 &= \underset{k \in \mathbb{F}_2^n}{\text{argmin}} \text{tr} \left((Q-1) \hat{\Sigma}^{-1} \hat{\Sigma} \right) = \underset{k \in \mathbb{F}_2^n}{\text{argmin}} D(Q-1). \tag{14}
 \end{aligned}$$

Indeed, at line (13), we demonstrated in the proof of Theorem 3 in that the minimal value (9) of α is independent on Σ . Eventually, it can be seen at line (14) that the distinguisher with $\hat{\Sigma}$ instead of Σ does not depend on the key⁶.

⁶ Indeed, $\text{argmin}_k \text{cst} = \mathbb{F}_2^n$, meaning that all keys are equiprobable. Intuitively, when both the noise and the model parameters are regressed at the same time, any key achieves the same match between parametric model and side-channel observations.

4.2 Attacks on Simulated Traces

In this subsection we present simulations when α is known exactly or regressed online. We consider an attack of PRESENT, where the SBox is $n = 4 \rightarrow n = 4$. For the sake of the simulations, we choose two kinds of α :

- “identical”: all the $n = 4$ bits leak the same waveform, like in the Hamming weight model,
- “proportional”: the waveform has weight 1 for SBox bit 0, and is multiplied by 2 (resp. 3 and 4) for SBox bit 1 (resp. 2 and 3).

The waveform for each bit is that represented in Fig. 1 (upper left graph). Specifically, for all $1 \leq d \leq D$ and $1 \leq s \leq S$, the envelope consists in damped oscillations:

$$\alpha_{d,s} = e^{-\frac{2d}{D}} \cos\left(2\pi \frac{d}{D}\right) \quad \text{for the “identical” case,} \quad (15)$$

$$\alpha_{d,s} = s \cdot e^{-\frac{2d}{D}} \cos\left(2\pi \frac{d}{D}\right) \quad \text{for the “proportional” case.} \quad (16)$$

The noise is chosen normal, using two distributions:

- “isotropic”: the covariance matrix is σ^2 times the $D \times D$ identity,
- “auto-regressive” (of “AR” for short): the covariance matrix element at position (d, d') , for $1 \leq d, d' \leq D$, is $\sigma^2 \rho^{|d-d'|}$. This noise is not independent from sample to its neighbours, but the correlation ρ decreases exponentially as samples get further apart.

Proposition 7. *The success probability of \mathcal{D}_{ML} is greater than that of $\mathcal{D}_{\text{ML,sto}}$.*

Proof. Indeed, according to [8], \mathcal{D}_{ML} maximizes the success probability. Thus, the distinguisher $\mathcal{D}_{\text{ML,sto}}$ has a smaller success probability. The success probability is the same if the minimization over α in the proof of Theorem 3 yields the exact matrix α used in the model (1). \square

Simulations allow to estimate the loss in terms of efficiency of not knowing the model (Proposition 7), by comparing distinguishers \mathcal{D}_{ML} ((4)) and $\mathcal{D}_{\text{ML,sto}}$ ((8)). The success rate of the optimal distinguisher \mathcal{D}_{ML} is drawn in order to materialize the limit between feasible (below) and unfeasible (above) attacks.

Results for low noise ($\sigma = 1$) are represented in Fig. 5. We can see that the Hamming weight model is clearly harder to attack, because the leakage of one bit cannot be distinguished from that of the other bits. Besides, we notice that the stochastic attack is performing much worse than the optimal attack: about 10 times more traces are required for an equivalent success probability in key extraction.

Results for high noise ($\sigma = 4$) are represented in Fig. 6. Again, the “proportional” model is easier to attack than the “identical” model (for each bit). Now, we also see that the gap between the optimal ML attack and the stochastic attack narrows: only about 5 times more traces are needed for the stochastic attack

to perform as well as the optimal attack in terms of success probability. Besides, we notice that the AR noise is favorable to the attacker. It is therefore important in practice for the attacker to characterize precisely the noise distribution (recall the methodology presented in Sec. 4.1).

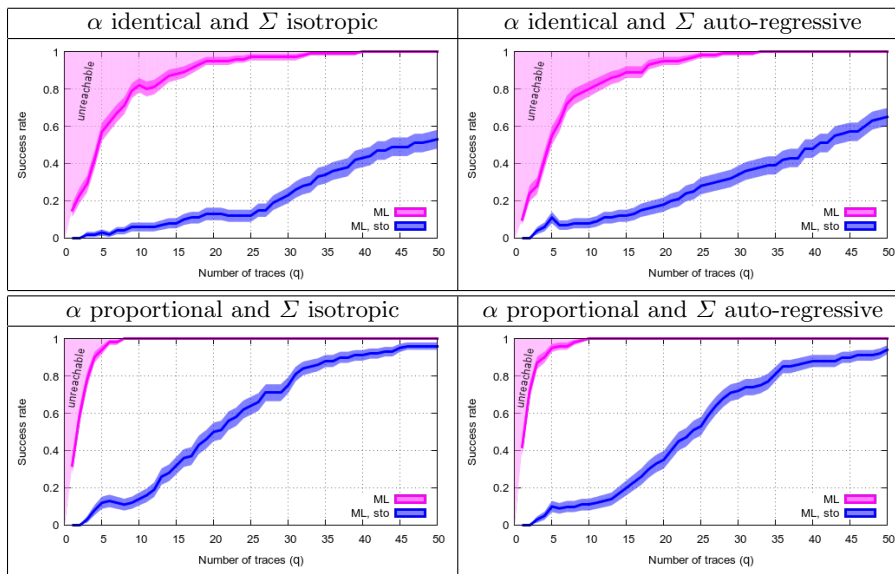


Fig. 5. Simulations for $D = 3$, $S = 5$, $n = 4$, $\sigma = 1$ (AR noise with $\rho = 0.5$).

Clearly, these conclusions are in line with the *template* versus *stochastic* (of-line) study carried out in [6]: for high noise, the (online) learning of the model requires more traces, hence is more accurate. Therefore, the performance of $\mathcal{D}_{ML,sto}$ gets closer to that of \mathcal{D}_{ML} than for low noise.

4.3 Attacks on Real-World Traces

We now compare CPA with \mathcal{D}_{ML} and $\mathcal{D}_{ML,sto}$ on measurements provided by the DPA contest V4. These traces have been acquired from an 8-bit processor, hence have a signal-to-noise ratio greater than one, reaching 7 at some points in time. The interval for our case-study is $[170, 250]$ from Fig. 2, hence $D = 80$. Regarding ML, two learning strategies have been implemented:

1. the model is learned from a disjoint set of 5k traces, which is the operational scenario for a profiled attack;
2. the model is learned from the traces being attacked (denoted **self** in Fig. 7). This case does not represent a realistic attack, but is interesting in that it highlights the best possible attacker.

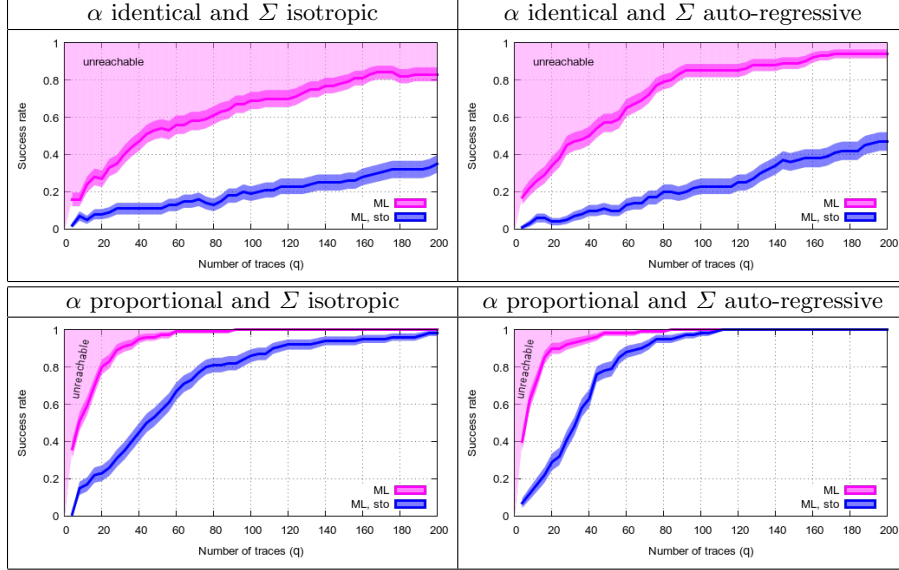


Fig. 6. Simulations for $D = 3$, $S = 5$, $n = 4$, $\sigma = 4$ (AR noise with $\rho = 0.5$).

The attack success rates are plotted in Fig. 7. One can see that both variants of \mathcal{D}_{ML} and $\mathcal{D}_{ML,sto}$ achieve better with $S = 9$ than with $S = 2$. This is consistent with the analysis carried out in Sec. 2.3. Actually, the CPA has a very poor performance because the model is actually very far from a Hamming weight: as can be seen in Fig. 2.3(a), some parameters α_i (e.g., for $i = 2$ and 6) are positive in region $[180, 200]$ whereas others α_j (e.g., for $j = 1, 3, 4$ and 5) are negative. The compensating signs account why the Hamming weight model is inappropriate. The ML with model pre-characterization on the traces under attack show that very strong attacks are possible (using a few traces only). Interestingly, when the model used by ML is characterized on 5k traces distinct from the traces being attacked, the performance is almost similar. Eventually, the online stochastic attack derived in this paper ($\mathcal{D}_{ML,sto}$) performs better than CPA (the distinguisher being the maximum value of the Pearson correlation over the $D = 80$ samples).

5 Conclusions and Perspectives

Distinguishing a key from both multivariate leakage samples and multiple models can be done in one step as shown in this paper. A compact expression of the distinguisher is provided, using matrix operations. The strategy is applied to real-world traces in profiled and non-profiled scenarios. The resulting attack is more efficient than the traditional approach “dimensionality reduction then stochastic (linear regression) attack”. The new multivariate distinguisher out-

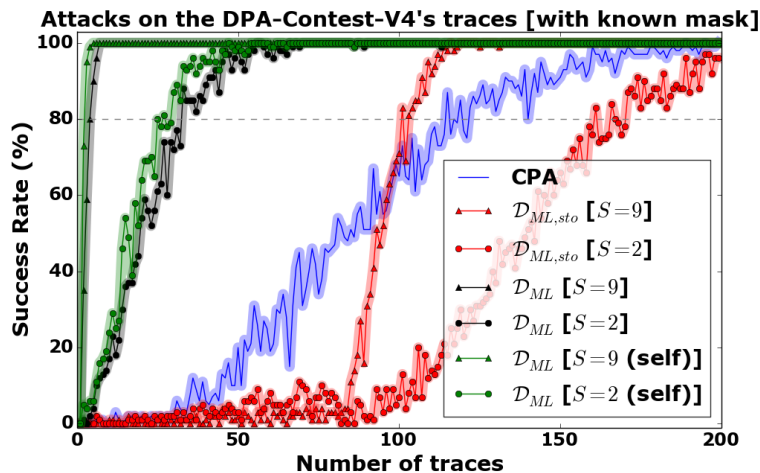


Fig. 7. Comparison of success rate of CPA, $\mathcal{D}_{ML,sto}$ for $S \in \{9, 2\}$, and \mathcal{D}_{ML} for $S \in \{9, 2\}$ (with two distinct learning methods)

performs the other state-of-the-art attacks. The presented methodology allows for leakage agnostic attacks on vectorial leakage measurements and complex models. In addition, the matrix-based expression of the distinguisher benefits from matrix-oriented software that implements computational optimizations for large dimensions.

A companion future work would consist in determining the optimal model dimensionality and basis from any acquisition campaign. Another perspective is to adapt the methodology to masked implementations, as already done for monivariate leakage in [4], yet for this case the distinguishers will certainly not exhibit simple closed-form expressions. However, we believe that the approach could be fruitful in practice backed with suitable optimization software.

References

1. Éric 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.
2. Nicolas Bruneau, Sylvain Guilley, Annelie Heuser, Damien Marion, and Olivier Rioul. Less is More - Dimensionality Reduction from a Theoretical Perspective. In Tim Güneysu and Helena Handschuh, editors, *Cryptographic Hardware and Embedded Systems - CHES 2015 - 17th International Workshop, Saint-Malo, France, September 13-16, 2015, Proceedings*, volume 9293 of *Lecture Notes in Computer Science*, pages 22–41. Springer, 2015.

3. Suresh Chari, Josyula R. Rao, and Pankaj Rohatgi. Template Attacks. In *CHES*, volume 2523 of *LNCS*, pages 13–28. Springer, August 2002. San Francisco Bay (Redwood City), USA.
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. Julien Doget, Emmanuel Prouff, Matthieu Rivain, and François-Xavier Standaert. Univariate side channel attacks and leakage modeling. *J. Cryptographic Engineering*, 1(2):123–144, 2011.
6. Benedikt Gierlichs, Kerstin Lemke-Rust, and Christof Paar. Templates vs. Stochastic Methods. In *CHES*, volume 4249 of *LNCS*, pages 15–29. Springer, October 10-13 2006. Yokohama, Japan.
7. Cezary Glowacz, Vincent Grosso, Romain Poussier, Joachim Schüth, and François-Xavier Standaert. Simpler and More Efficient Rank Estimation for Side-Channel Security Assessment. In Gregor Leander, editor, *Fast Software Encryption - 22nd International Workshop, FSE 2015, Istanbul, Turkey, March 8-11, 2015, Revised Selected Papers*, volume 9054 of *Lecture Notes in Computer Science*, pages 117–129. Springer, 2015.
8. 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.
9. Victor Lomné, Emmanuel Prouff, and Thomas Roche. Behind the Scene of Side Channel Attacks. In Kazue Sako and Palash Sarkar, editors, *ASIACRYPT (1)*, volume 8269 of *Lecture Notes in Computer Science*, pages 506–525. Springer, 2013.
10. Luke Mather, Elisabeth Oswald, and Carolyn Whittall. Multi-target DPA attacks: Pushing DPA beyond the limits of a desktop computer. In Sarkar and Iwata [11], pages 243–261.
11. 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*. Springer, 2014.
12. Werner Schindler, Kerstin Lemke, and Christof Paar. A Stochastic Model for Differential Side Channel Cryptanalysis. In *LNCS*, editor, *CHES*, volume 3659 of *LNCS*, pages 30–46. Springer, Sept 2005. Edinburgh, Scotland, UK.
13. Takeshi Sugawara, Naofumi Homma, Takafumi Aoki, and Akashi Satoh. Profiling attack using multivariate regression analysis. *IEICE Electronics Express*, 7(15):1139–1144, 2010.
14. TELECOM ParisTech SEN research group. DPA Contest (4th edition), 2013–2014. <http://www.DPAcontest.org/v4/>.
15. Nicolas Veyrat-Charvillon, Benoît Gérard, Mathieu Renaud, and François-Xavier Standaert. An Optimal Key Enumeration Algorithm and Its Application to Side-Channel Attacks. In Lars R. Knudsen and Huapeng Wu, editors, *Selected Areas in Cryptography*, volume 7707 of *Lecture Notes in Computer Science*, pages 390–406. Springer, 2012.
16. Nicolas Veyrat-Charvillon, Benoît Gérard, and François-Xavier Standaert. Soft analytical side-channel attacks. In Sarkar and Iwata [11], pages 282–296.