

# Side-Channel Attacks Against the Human Brain: the PIN Code Case Study

Joseph Lange, Clément Massart, André Mouraux, Francois-Xavier Standaert  
Université catholique de Louvain, B1348 Louvain-la-Neuve, Belgium.  
e-mails: clement.massart, andre.mouraux, fstandae@uclouvain.be

**Abstract.** We revisit the side-channel attacks with Brain-Computer Interfaces (BCIs) first put forward by Martinovic et al. at the USENIX 2012 Security Symposium. For this purpose, we propose a comprehensive investigation of concrete adversaries trying to extract a PIN code from electroencephalogram signals. Overall, our results confirm the possibility of partial PIN recovery with high probability of success in a more quantified manner (i.e., entropy reductions), and put forward the challenges of full PIN recovery. They also highlight that the attack complexities can significantly vary in function of the adversarial capabilities (e.g., supervised / profiled vs. unsupervised / non-profiled), hence leading to an interesting tradeoff between their efficiency and practical relevance. We then show that similar attack techniques can be used to threaten the privacy of BCI users. We finally use our experiments to discuss the impact of such attacks for the security and privacy of BCI applications at large, and the important emerging societal challenges they raise.

## 1 Introduction

**State-of-the-art.** The increasing deployment of Brain Computer Interfaces (BCIs) allowing to control devices based on cerebral activity has been a permanent trend over the last decade. While originally specialized to the medical domain (e.g., [13,22]), such interfaces can now be found in a variety of applications. Notorious examples include drowsiness estimation for safety driving [19] and gaming [9]. Quite naturally, these new capabilities come with new security and privacy issues, since the signals BCIs exploit can generally be used to extract various types of sensitive information [7,15]. For example, at the USENIX 2012 Security Symposium, Martinovic et al. showed empirical evidence that electroencephalogram (EEG) signals can be exploited in simple, yet effective attacks to (partially) extract private information such as credit card numbers, PIN codes, dates of birth and locations of residence from users [21]. These impressive results leveraged a broad literature in neuroscience, which established the possibility to extract such private information (e.g., see [14] for lie detection and [16] for neural markers of religious convictions). Or less invasively, they can be connected to linguistic research on the reactions of the brain to semantic associations and incongruities (e.g., [17,18,6]). All these threats gain concrete relevance with the availability of EEG-based gaming devices to the general public [1,2].

**Motivation & goals.** Based on this state-of-the-art, the next step is to push the evaluation of the side-channel threat model in the context of BCI-based applications further. In this respect, the seminal work of Martinovic et al. clearly puts forward the existence of an exploitable bias for various types of private information extraction. But quantifying the impact of this bias in advanced adversarial contexts was left as an important challenge. Typical questions include:

- Can we exactly extract private information with high success rate by increasing the number of observations in side-channel attacks exploiting BCIs?
- How does the effectiveness of unsupervised (aka non-profiled) side-channel attacks exploiting BCIs compare to supervised (aka profiled) ones?
- How efficiently can an adversary build a sufficiently accurate model for supervised (aka profiled) side-channel attacks exploiting BCIs?

Interestingly, these are typically questions that have been intensively studied in the context of side-channel attacks against cryptographic devices (see [20] for an engineering survey and the proceedings of the CHES conference for regular advances in the field [3]). In particular, a recurring problem in the analysis of such implementations is to determine their worst-case security level, in order to bound the probability of success of any adversary in the most accurate manner [27]. This implies very different challenges than in the standard cryptographic setting, since the efficiency of such physical attacks highly depends on the adversary’s understanding and knowledge of his target device. Hence, a variety of tools have been developed in order to ensure that side-channel security evaluations are “good enough” (as described next). Our goal in this paper is to investigate the applicability of such tools in order to answer the previous questions regarding the efficiency and impact of side-channel attacks against the human brain.

**Contributions.** For this purpose, we propose an in-depth study of (a variation of) one of the case studies in [21], namely side-channel PIN code recovery attacks, that share some similarities with key recovery attacks against embedded devices. In this respect, our contributions are threefold. After a description of our experimental settings (Section 2), we first describe a methodology allowing us to analyze the informativeness of EEG signals and their impact on security with confidence (Section 3). While this methodology indeed borrows tools from the field of side-channel attacks against cryptographic implementations, it also deals with new constraints (e.g., the limited amount of observations available for the evaluations, and the less regular distribution of these observations, for which a very systematic and principled approach is particularly important). Second, we provide a comprehensive experimental evaluation of our side-channel attacks against the human brain using this methodology (Section 4). We combine information theoretic and security analyzes in the supervised / profiled and unsupervised / non-profiled contexts, provide quantified estimates for the complexity of the attacks, and pay a particular attention to the stability of and confidence in our results. We conclude by discussing consequences the consequences of our work for the security and privacy of BCI-based applications Section 5).

Admittedly, and as will be discussed in detail next, our results can be seen as positive or negative. That is, we show in the same time that partial information about PINs can be extracted with confidence, and that full PIN extractions are challenging because of the high cardinality of the target and risks of false positive. So they should mostly be viewed as a warning flag that such partial information is possible and may become critical when the cardinality of the target decreases and/or large amounts of data are available to the adversary.<sup>1</sup>

## 2 Experimental setting, threat model and limitations

In our experiments, eight people (next denoted as users) agreed to provide the 4-digit PIN code that they consider the most significant to them, meaning the one they use the most frequently in their daily life. This PIN code was given by the users before the experiment started, stored during the experiment, and deleted afterward for confidentiality reasons. Five other random 4-digit codes were generated for each user (meaning a total of six 4-digit codes per user).

Each (real or random) PIN was then shown on a computer exactly 150 times to each user (in a random order), meaning a total of 900 events for which we recorded the EEG signal in sets of 300, together with a tag  $T$  ranging from 1 to 6 (with  $T = 1$  the correct PIN and  $T = 2$  to 6 the incorrect ones). We used 32 Ag-AgCl electrodes for the EEG signals collection. These were placed on the scalp using a Waveguard cap from Cephalon, using the international 10-10 system. The Stimulus Onset Asynchrony (SOA) was set to 1,009s (i.e., slightly more than one second, to reduce the environmental noise). The time each PIN was shown was set to 0,5s. When no PIN was displayed on the screen, a + sign was maintained in order to keep the focus of the user on the center of the screen. We additionally ensured that two identical 4-digit codes were always separated by at least two other 4-digit codes. The split of our experiments in sub-experiments of 300 events was motivated by a maximum duration of 5 minutes, during which we assumed the users to remain focused on the screen. The signals were amplified and sampled at a 1000Hz rate with a 32-channel ASA-LAB EEG system from Advanced Neuro Technologies. Eventually, and in order to identify eye-blinks which potentially perturb the EEG signal, we added two bipolar surface electrodes on the upper left and lower right sides of the right eye, and rejected the records for which such an artifact was observed. This slightly reduced the total number of events stored for each user (precisely, this number was reduced to 900, 818, 853, 870, 892, 887, 878, 884, for users 1 to 8).

This simplified setting naturally comes with limitations. First and concretely, the number of possible PIN codes for a typical smart card would of course be much larger than the 6 ones we investigate (e.g., 10,000 for a 4-digit PIN). In this respect, we first insist that the primary goal of the following experiments

---

<sup>1</sup> The experiments described next were approved by the local Research Ethics Committee and performed in compliance with the Code of Ethics of the World Medical Association (Declaration of Helsinki). All participants gave written informed consent.

is to investigate the information leakages in EEG signals thoroughly, and this limited number of PIN codes allowed us to draw conclusions with good statistical confidence. Yet, we also note that this setting could be extended to a reasonable threat model. For example, one could target  $\approx 1000$  different users by repeatedly showing them  $\approx 10$  PIN codes among the 10,000 possible ones, and recover one PIN with good confidence. Second, and since the attacks we carry out essentially test familiar vs. unfamiliar information, there is also a risk of false positives (e.g., an all zero code or a close to correct code). This is in fact something we observed in our experiments. In this respect, our mitigation plan is to exploit statistical tools minimizing the number of false negatives, therefore potentially allowing enumeration among the most likely candidates [28].

### 3 Methodology

In this section, we describe the methodology we used in order to assess and better quantify the feasibility of side-channel attacks against the human brain. Concretely, and contrary to the case of embedded devices where the leakage distributions are supposed to be stable and the number of observations made by the adversary can be large, we deal with a very different challenge. Namely, we need to cope with irregular distributions possibly affected by outliers, and can only assume a limited number of observations.

As a result, the following sections mainly aim to convince the reader that our treatment of the EEG signals is not biased by dataset-specific overfitting. For this purpose, our strategy is twofold. First, we apply the same (pre)processing methods to the measurements of all the users. This means the same selection of electrodes, the same dimensionality reduction and Probability Density Function (PDF) estimation tools (with identical parameters), and the same outliers definition. Second, we systematically verified that our results were in the same time consistent with neurophysiological expectations, and stable across a sufficient range of (pre)processing parameters. As a result, our primary focus is on the confidence in and stability of the results, more than on their optimality (which is an interesting scope for further research). In other words, we want to guarantee that EEG signals provide exploitable side-channel information for PIN code recovery, and to evaluate a sufficient number of observations for which such an attack can be performed with good success probability.

#### 3.1 Notations

We denote the (multivariate) EEG signals of our experiments with a random variable  $\mathbf{O}$ , a sample EEG signal as  $\mathbf{o}$ , and the set of all the observations available for evaluation as  $\mathcal{O}$ . These observations depend on (at least) three parameters: the user under investigation, next denoted with a random variable  $U$  such that  $u \in \{1, 2, \dots, 8\}$ ; the nature of the 4-digit code observed (i.e., whether it is correct or a random PIN), next denoted with a random variable  $P$  such that  $p \in \{0, 1\}$ ; and a noise random variable  $N$ . Each observation is initially made of 32 vectors of 1,000 samples, corresponding to 32 electrodes and  $\approx 1$ s per event.

### 3.2 Supervised (aka profiled) evaluation

In order to best evaluate the actual informativeness of the EEG signals regarding the PIN displayed in our experiments, and inspired by the worst-case side-channel security evaluations of cryptographic devices, our work first investigates so-called profiled attacks, which correspond to a supervised machine learning context. For this purpose, a part of the observations in  $\mathcal{O}$  are used to estimate a (probabilistic) model  $\Pr_{\text{model}}[P = p | \mathbf{O} = \mathbf{o}]$ . The adversary/evaluator then uses this model in order to try extracting the PIN from the remaining observations. Note that our profiling is based on the binary random variable  $p$ , where  $p = 0$  if the PIN is random and  $p = 1$  if the PIN is real, and not based on the value of the PIN tag itself. This is motivated by the following practical and neurophysiological reasons:

- From a practical point-of-view, building a model for all the PINs and users seems impractical in real-world settings: this would require being able to collect multiple observations for each of the 10,000 possible values of a 4-digit code. Furthermore, and as discussed in Section 3.3, our real vs. random profiling allowed us to lean towards realistic (non-profiled) attacks.
- From a neurophysiological point-of-view, the information we aim to extract is based on Event-Related Potentials (ERPs) that have been shown to reflect semantic associations and incongruities [17,18,6]. In this respect, while we can expect a user to react differently to real and random 4-digit codes, there is no reason for him to treat the random codes differently.

**A. Evaluation metrics** Following the general principles put forward in [27], our evaluations will be based on a combination of information theoretic and security analyzes. The first ones aim at evaluating whether exploitable information is available in the EEG signals; the second ones at evaluating how efficiently this information can be exploited to mount a side-channel attack. Note that since we do not assume the users to behave identically, these metrics will always be evaluated and discussed for each user independently.

**Perceived information.** The Perceived Information (PI) was introduced in the context of side-channel attacks against cryptographic devices, of which the goal is to recover some secret data (aka key) given some physical leakage [23]. The PI aims at quantifying the amount of information about the secret key, independent of the adversary who will exploit this information. Informally, we will use this metric in a similar way, by just considering  $P$  as a bit to recover, and the observations as leakages. Using the previous notations, we define the PI between the PIN random variable  $P$  and the observation random variable  $\mathbf{O}$ :

$$\text{PI}(P; \mathbf{O}) = H[P] + \sum_p \Pr[p] \cdot \int_{\mathbf{o}} f(\mathbf{o}|p) \cdot \log_2 \Pr_{\text{model}}[p|\mathbf{o}] d\mathbf{o},$$

where we use the notation  $\Pr[X = x] =: \Pr[x]$  for conciseness, and  $f(\mathbf{o}|p)$  is the (continuous) PDF of the observations given the value of  $p$ . In the ideal case

where the model is perfect, the PI is identical to Shannon’s mutual information. In the practical cases where the model differs from the observation’s true distribution, the PI captures the amount of information that is extracted from these observations, biased by the model (assumption & estimation) errors [11].

Of course, concretely the true distribution  $f(\mathbf{o}|p)$  is unknown to the adversary/evaluator and can only be sampled. Therefore, the approach in side-channel analysis, that we repeat here, is to split the set of observations  $\mathcal{O}$  in  $k$  non-overlapping sets  $\mathcal{O}^{(i)}$ . We then define the profiling sets  $\mathcal{O}_p^{(j)} = \bigcup_{i \neq j} \mathcal{O}^{(i)}$  and the test sets  $\mathcal{O}_t^{(j)} = \mathcal{O} \setminus \mathcal{O}_p^{(j)}$ . The PI is computed in two phases:

1. The observations’ conditional distribution is estimated from a profiling set. We denote this phase with  $\hat{f}_{\text{model}}^{(j)}(\mathbf{o}|p) \leftarrow \mathcal{O}_p^{(j)}$ . Note that the  $\Pr_{\text{model}}[p|\mathbf{o}]$  factor involved in the PI definition is directly derived via Bayes’ theorem as:

$$\hat{\Pr}_{\text{model}}[p|\mathbf{o}] = \frac{\hat{f}_{\text{model}}^{(j)}(\mathbf{o}|p) \cdot \Pr[p]}{\sum_{p^*} \hat{f}_{\text{model}}^{(j)}(\mathbf{o}|p^*) \cdot \Pr[p^*]}.$$

2. The model is then tested by computing the PI estimate:

$$\hat{\text{PI}}^{(j)}(P; \mathbf{O}) = H[P] + \sum_{p=0}^1 \Pr[p] \cdot \sum_{\mathbf{o} \in \mathcal{O}_t^{(j)}|p} \frac{1}{n_p^j} \cdot \log_2 \hat{\Pr}_{\text{model}}[p|\mathbf{o}],$$

where  $n_p^j$  is the number of observations in the test set  $\mathcal{O}_t^{(j)}|p$ .

Eventually, the  $k$  outputs  $\hat{\text{PI}}^{(j)}(P; \mathbf{O})$  are averaged to get an unbiased estimate, and their spread characterizes the accuracy of the result (see Paragraph G). Note that concretely, the maximum size for the profiling set in our experiments equals  $\approx 899$ , leading to a cross-validation parameter  $k \approx 900$  and a test set of size 1. In this case, the model building phase is repeated  $\approx 900$  times, and each model is tested once against an independent sample. (We use the  $\approx$  symbol to reflect the fact that these values are approximated, due to the rejection of eye blinks mentioned in Section 2). This “leave one out” strategy has a large cross-validation parameter compared to current practice (e.g., in side-channel attacks against cryptographic implementations a value of  $k = 10$  was selected [11]), leading to computationally intensive evaluations. Yet, it is justified in our study because of the limited number of samples available in our experiments.

**Success rate and average rank.** In order to confirm that the estimated PI indeed leads to concrete attacks, we consider two simple security metrics. Here, the main challenge is that we only have models for the real and random PIN codes, while the actual observations in the test set naturally come from six different events. As a result, we first considered the success rate event per event. For this purpose, the  $\approx 900$  observations are split in 6 sets of  $\approx 150$  observations that correspond to the six different tag values. Based on these 6 sets, we can compute the probability that the observations are correctly classified as real or

random in function of the number of observations exploited in the attack, next denoted as  $q$ . This is done by averaging a success function  $S$  that is computed as follows. If  $q = 1$ :  $S(\mathbf{o}_1) = 1$  if  $\hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1] > \hat{\text{Pr}}_{\text{model}}[\bar{p}|\mathbf{o}_1]$  and  $S(\mathbf{o}_1) = 0$  otherwise (where  $\bar{p}$  denotes the incorrect event); if  $q = 2$ :  $S(\mathbf{o}_1, \mathbf{o}_2) = 1$  if  $\hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1] \times \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_2] > \hat{\text{Pr}}_{\text{model}}[\bar{p}|\mathbf{o}_1] \times \hat{\text{Pr}}_{\text{model}}[\bar{p}|\mathbf{o}_2]$ ; ... Concretely, this success rate is an interesting metric to check whether the observations generated by different incorrect PIN values indeed behave similarly.

Of course, an adversary eventually wants to compare the likelihoods of different PIN values. For this purpose, we also considered the average rank of the correct PIN in an experiment where we gradually increase the number of observations per tag  $q$ , but this time consider sets of 6 observations at once, that we classify only according to the model for the real PIN. This leads to vectors  $(\hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^1], \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^2], \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^3], \dots, \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^6])$  if  $q = 1$ ,  $(\hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^1] \times \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_2^1], \dots, \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_1^6] \times \hat{\text{Pr}}_{\text{model}}[p|\mathbf{o}_2^6])$  if  $q = 2$ , ... , where the superscripts denote the tag from which the observations originate. The average rank is then obtained by sorting this vector and estimating the sample mean of the position of the tag 1 in the sorted vector.

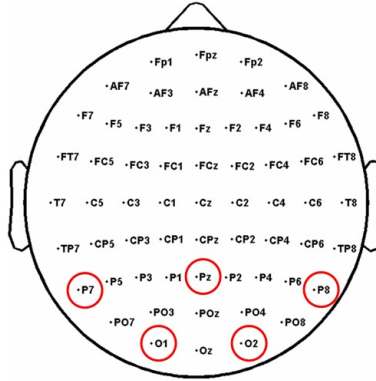
**Connecting the metrics (sanity check).** Note that as discussed in [10], information theoretic and security metrics can be connected (i.e., a model that leads to a positive PI should lead to successful attacks).<sup>2</sup> We consider both types of metrics in our experiments because the first ones allow a better assessment of the confidence in the evaluations (see Paragraph G) while the second ones lead to simpler intuitions regarding the concrete impact of the attacks.

**B. Preprocessing** As a first step, all the observations were preprocessed using a bandpass filter. We set the low-frequency cut-off to 0.5Hz to remove the slow drifts in the EEG signals, and the high-frequency cut-off to 30Hz to remove muscle artifacts and 50Hz environmental noise.

**C. Selection of electrodes** As mentioned in introduction, each original observation is made of 32 vectors of 1,000 samples, leading to a large amount of data to process. To simplify our treatments, we started by analyzing the different electrodes independently. Among the 32 ones of our cap, electrodes P7, P8, Pz, O1 and O2 gave rise to non-negligible signal (see Figure 1), which is consistent with the existing literature where ERPs related to semantic associations and incongruities were exhibited in the central/parietal zones [17,18,6]. Our following analyzes are based on the exploitation of the electrodes P7 and P8 which provided the most regular information across the different users.

For illustration, Figures 2 and 3 represent the mean and standard deviation traces corresponding to two different users. From these examples, a couple of relevant observations can already be extracted (and will be useful for the design

<sup>2</sup> More precisely, the PI is an average metric, so what is needed is that each line of the PI matrix defined in [27] (corresponding to 6 different events in our study) are positive, which we observed and confirmed with the success rate analysis.



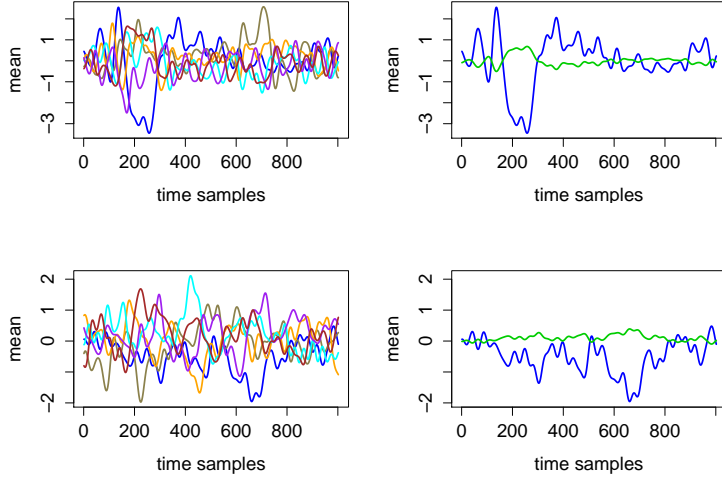
**Fig. 1.** Repartition of the electrodes on the scalp.

and interpretation of our following evaluations). First, we see (on the left parts of Figure 2) that the EEG signals may be more or less informative depending on the users and electrodes. More precisely, we generally noticed informative ERP components after 300 to 600 milliseconds (known as the P300) for most users and electrodes, which is again consistent with the existing literature [17,18,6]. Yet, our measurements also put forward user-specific differences in the shape of the mean traces corresponding to the correct PIN value. (Note that the figure only shows examples of informative EEG signals, but for some other users and electrodes, no such clear patterns appear). Second, and quite importantly, the difference between the left and right parts of the figures illustrates the significant gain when moving from an unsupervised / unprofiled evaluation context to a supervised / profiled one. That is, while in the first case, we need the traces corresponding to the correct PIN value to stand out, in the second case, we only need it to behave differently than the others.

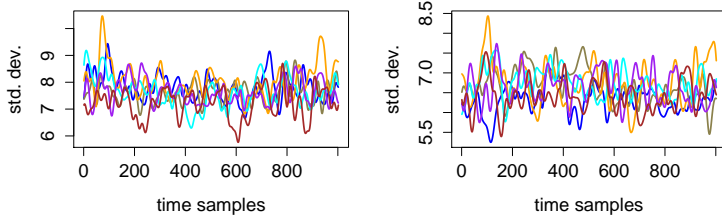
Eventually, a look at the standard deviation curves in Figure 3 suggests that the measurements are quite noisy, hence non-trivial to exploit with a limited amount of observations. This will be confirmed in our following PDF estimation phase, and therefore motivates the dimensionality reduction in the next section (intuitively because using more dimensions can possibly lead to better signal extraction, which can mitigate the effect of a large noise level).

**D. Dimensionality reduction** The evaluation of our metrics requires to build a probabilistic model, which may become data intensive as the number of dimensions in the observations increases. For example, directly estimating a 2000-dimensional PDF corresponding to our selected electrodes is not possible. In order to deal with this problem, we follow the standard approach of reducing dimensionality. More precisely, we use the Principal Component Analysis (PCA) that was shown to provide excellent results in the context of side-channel attacks against cryptographic devices [4]. We investigate two options in this direction.





**Fig. 2.** Exemplary mean traces for different tag (left) and PIN (right) values. Top: User 8, Electrode P7. Bottom: User 6, Electrode P7.

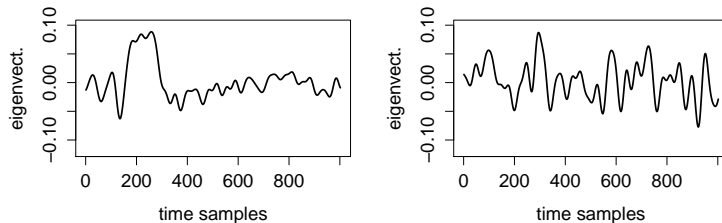


**Fig. 3.** Exemplary standard deviation traces for different tag values corresponding to User 8, Electrode P7 (left) and User 6, Electrode P7 (right).

First, and looking at the observations in Figure 2, it appears that the mean traces corresponding to the different tags are quite discriminant regarding the value of  $p$ . Hence, and as in [4], a natural option is to compute the projection vectors of the PCA based on these mean traces. This implies computing average vectors  $\bar{\sigma}^j = \mathbb{E}_{i \approx 1}^{150} \sigma_i^j$ , and then to derive the PCA eigenvectors based on the  $\bar{\sigma}^j$ 's, which we denote as  $\mathbf{R}_{1:N_d} \leftarrow \text{PCA}(\{\bar{\sigma}^j\}_{j=1:6})$ , where  $N_d$  is the number of dimensions to extract. Due to the limited number of mean traces (i.e., 6), we can only compute  $N_d = 5$  eigenvectors, and therefore are limited to 5-dimensional attacks in this case.<sup>3</sup> However, it turned out that in our experiments, this version of the PCA extracts most of the relevant samples in the first dimension. This is intuitively witnessed by Figure 4 which represents the first and fifth eigenvectors

<sup>3</sup> Because we used the small sample size variant of PCA in [4].

corresponding to User 8 and Electrode P7 (i.e.,  $\mathbf{R}_1$  and  $\mathbf{R}_5$ ): we indeed observe that the first dimension corresponds to the points of interest in Figure 2, while the fifth one seems to be dominated by noise. In the following, we will denote this solution as the “average PCA”. Note that such a dimensionality reduction does not take advantage of any secret information (i.e., it is not a supervised / profiled one) since it builds the mean traces based on public tags.



**Fig. 4.** Exemplary eigenvectors for the average PCA, corresponding to User 8, Electrode P7. Left: first dimension. Right: fifth dimension.

Yet, one possible drawback of the previous method is that estimating the average traces  $\bar{o}^j$  becomes expensive when the number of PIN codes increases. In order to deal with and quantify the impact of this limitation, we also considered a “raw PCA”, where we directly reduce the dimensionality based on raw traces, next denoted as  $\mathbf{R}_{1:N_d} \leftarrow \text{PCA}(\{\mathbf{o}_i\}_{i \approx 1:900})$ . While this approach is not expected to extract the information as effectively, it allows deriving a much larger number of dimensions than in the previous (average) case. Concretely though, exploiting dimensions 1 to 5 only was a good tradeoff between the informativeness of the dimensionality reduction, the risk of overfitting (useless) dataset-dependent patterns and the risk of outliers in our experiments (see Paragraph F).

As a result of this dimensionality reduction phase, the observation vectors  $\mathbf{o}(1:2000)$  (which correspond to the concatenation of the measurements for our two selected electrodes) are reduced to smaller vectors  $\mathbf{R}_{1:N_d} \times \mathbf{o}$  (i.e., each dimension  $o(d)$  corresponds to the scalar product between the original observations  $\mathbf{o}$  and a 2000-element vector  $\mathbf{R}_d$ ). We recall that PCA is not claimed to be an optimal dimensionality reduction, since it optimizes a criteria (i.e., the variance between the raw or mean traces) which does not capture all the information in our measurements. However, it is a natural first step in our investigations, and we could verify that our following conclusions are not affected by slight variations of the number of extracted dimensions (i.e., adding one or two dimensions), which therefore fits our (primary) confidence and stability goal.

**E. PDF estimation** We now describe the main ingredient of our supervised / profiled evaluation, namely the PDF estimation for which we exploit the knowledge of the  $p$  values for the observations in the profiling sets.

In order to build a model  $\hat{f}_{\text{model}}(\mathbf{o}_{1:N_d}|p)$ , we first take advantage of the fact that the dimensions of the  $\mathbf{o}_{1:N_d}$  vectors after PCA are orthogonal. By additionally considering them as independent, this allows us to reduce the PDF estimation problem from one  $N_d$ -variate one to  $N_d$  univariate ones. Based on this simplification, the standard approach in side-channel analysis is to assume the observations to be normally distributed, and to build Gaussian templates [8]. Yet, in our experiments no such obvious assumption on the distributions in hand was a priori available. As a result, we first considered a (non-parametric) kernel density estimation as used in [5], which has slower convergence but avoids any risk of biased evaluations [11]. Kernel density estimation is a generalization of histograms. Instead of bundling samples together in bins, it adds (for each observation) a small kernel centered on the value of the observation to the estimated PDF. The resulting estimation that is a sum of kernels is smoother than histograms and usually converges faster. Concretely, kernel density estimation requires selecting a kernel function (we used a Gaussian one) and to set the bandwidth parameter (which can be seen as a counterpart to the bin size in histograms). The optimal choice of the bandwidth depends on the distribution of the observations, which is unknown in our case. So we need to rely on a heuristic, and used Silverman’s rule-of-thumb for this purpose [24].

**F. Outliers** As mentioned in Paragraph D, the main drawback of the raw PCA is that it extracts the useful EEG information less efficiently, which we mitigate by using more dimensions. Unfortunately, this comes with an additional caveat. Namely, the less informative information extraction combined with the addition of more dimensions increases the risk of outliers (i.e., observations that would classify the correct PIN value very badly for some dimensions, possibly leading to a negative PI). In this particular case, we considered an additional post-processing (after the dimensionality reduction and model building phases). Namely, given the  $\approx 900$  probabilities  $\hat{\text{Pr}}[p|\mathbf{R}_{1:N_d} \times \mathbf{o}_i]$ , we rejected the ones below 0.001 and beyond 0.999. This choice is admittedly heuristic, yet did consistently lead to positive results for all the users. It is motivated by limiting the weight of the log probabilities for the outliers in the PI estimation. We insist that this treatment of outliers is only needed for the raw PCA. For the average PCA, we did not reject any observation (other than the ones in Section 2).

**G. Confidence** By using  $\approx 900$ -fold cross-validation, we can guarantee that our PI estimates will be based on 900 observations, leading to 900 values for the log probabilities  $\log_2(\hat{\text{Pr}}[p|\mathbf{R}_{1:N_d} \times \mathbf{o}_i])$ . Since this remains a limited amount of data compared to the case of side-channel attacks against cryptographic implementations, and the extracted PI values are small, we completed our information theoretic evaluations by computing a confidence interval for the PI estimates. To avoid any distribution-specific assumption, we computed a 10% bootstrap confidence interval [12], by resampling 100 bootstrap samples out of our 900 log probabilities, computing 100 mean bootstrap samples, sorting them, and using the 95th and 5th percentiles as the endpoints of the intervals. For simplicity, this was only done for the PI metric and not for the success rate and average

rank since (i) successful Bayesian attacks are implied by the information theoretic analysis [10], (ii) these metrics are more expensive to sample (e.g., we have only one evaluation of the success function with  $q \approx 150$  per user), and (iii) they are only exhibited to provide intuitions regarding the exploitability of the observations (i.e., the attack complexities).

### 3.3 Unsupervised (aka non-profiled) analysis

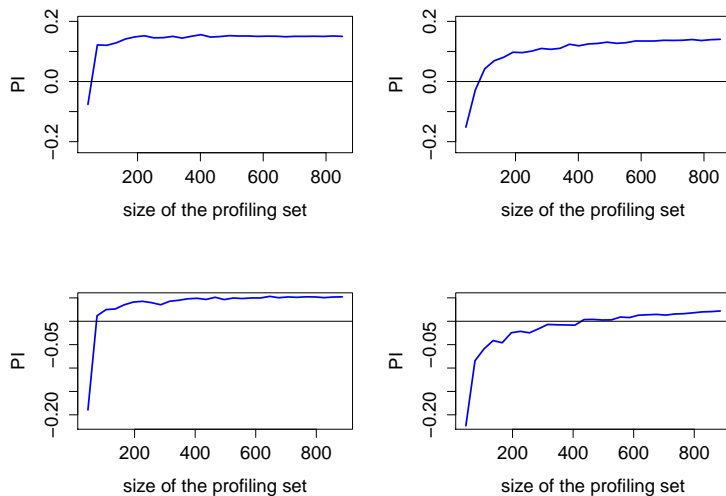
While supervised (aka profiled) analyzes are the method of choice to gain understanding about the information available in a side-channel, their practical applicability is of course questionable. Indeed, building a model for a target user may not always be feasible, and this is particularly true in the context of attacks against the human brain since (as discussed the long version of this paper), models built for one user are not always (directly) exploitable against another user. In this section, we therefore propose an unsupervised / non-profiled extension of the information theoretic evaluation outlined in Section 3.2. To the best of our knowledge, this variation was never described as such in the open literature (although it shares some similarities with the non-profiled attacks surveyed in [5]). For this purpose, our starting point is the observation from Figure 2, that in an unsupervised / non-profiled context, one can take advantage of the fact that the (e.g., mean) traces of the EEG signals corresponding to the correct PIN value may stand out. As a result, a natural idea is to compute the PI metric 6 times independently, each time assuming a different (possibly random) tag to be correct during an “on-the-fly” modeling phase. If the traces corresponding to the (truly) correct PIN are more singular (comparatively to the others), we can expect the PI estimated with this PIN to be larger, leading to a successful attack.

Of course, such an attack implies an additional neurophysiological assumption (while in the supervised / profiled setting, we just exploit any information available). Yet, it nicely fits the intuitions discussed in the rest of this section, which makes it a good candidate for concrete evaluation. Furthermore, we mention that directly recovering the correct PIN value may not always be necessary: as in the case of side-channel analysis, reducing the rank of the correct PIN value down to an enumerable one may be sufficient [28].

## 4 Experimental results

### 4.1 Supervised (aka profiled) evaluation

As in the previous section, we start with the results of our supervised / profiled evaluations, which will be in two (information theoretic and security) parts. Beforehand, there is one last choice regarding the computation of  $\hat{\Pr}[p|\mathbf{R}_{1:N_d} \times \mathbf{o}_i]$  via Bayes’ theorem described in Section 3.2, Paragraph A. Namely, should we consider maximum likelihood or maximum a posteriori attacks (i.e., should we take advantage of the a priori knowledge of  $\Pr[p]$  or consider a uniform a priori). Interestingly, in our context ignoring this a priori and performing maximum likelihood attacks is more relevant, since we mostly want to avoid false negatives



**Fig. 5.** Evolution of the PI in function of the size of the profiling set for Users 3 (top) and 6 (bottom), using average PCA (left) and raw PCA (right).

(i.e., correct PINs that would be classified as random ones), which prevent efficient enumeration. Since the a priori on  $P$  increases the amount of such errors (due to the a priori bias of 5/6 towards random PIN values), the rest of this section reports on the results of maximum likelihood attacks.

**A. Perceived Information** As a first step in our evaluations, we estimated the PI using the methodology described in the previous section. We started by looking at the evolution of the PI estimation in function of the number of observations in the profiling set used to build the model. The results of this analysis are in Figure 5 from which two quantities must be observed:

- The value of the PI estimate using the maximum profiling set (i.e., the extreme right values in the graphs). It reflects the informativeness of the model built in the profiling phases, and is correlated with the success rate of the online (maximum likelihood) attack using this model [10]. Positive PI values indicate that the model is sound (up to Footnote 2) and should lead to successful online attacks if the number of observations (i.e., the  $q$  parameter in our notations of Section 3.2) used by the adversary is sufficient.
- The number of traces in the profiling set required to reach a positive PI. It reflects the (offline) complexity of the model estimation (profiling) phase [26].

In this respect, the results in Figure 5 show a positive convergence for the two illustrated users, yet towards different PI values which indicates that the informativeness of the EEG signals differs between them. Next, and quite interestingly, we also see that the difference between average PCA (in the left part

of the figure) and raw PCA (in the right side) confirms the expected intuitions. Namely, the fact that raw PCA reduces dimensionality based on a less meaningful criteria and requires more dimensions implies a slower model convergence. Typically, model convergence was observed in the 100 observations’ range with average PCA and required up to 400 traces with raw PCA. For completeness, Table 1 contains the estimated PI values with maximum profiling set, for the different users and types of PCA. Excepted for one user (User 5) for which we could never reach a positive PI value,<sup>4</sup> this analysis suggests that all the users lead to exploitable information and confirms the advantage of average PCA.

User	$\hat{\text{PI}}(P; O)$ with avg. PCA	$\hat{\text{PI}}(P; O)$ with raw PCA
1	0.0739	0.0618
2	0.1643	0.1315
3	0.1494	0.1398
4	0.0920	0.0228
5	$\emptyset$	$\emptyset$
6	0.0521	0.0214
7	0.0759	0.0568
8	0.1697	0.0458

**Table 1.** Estimated PI values with maximum profiling set.

**B. Success rate and average rank** As discussed in Section 3.2, our information theoretic analysis is a method of choice to determine whether discriminant information can be extracted from EEG signals with confidence. Yet, it does not lead to obvious intuitions regarding the actual complexity of an online attack where an adversary obtains a set of  $q$  fresh observations and tries to detect whether some of them correspond to a real PIN value. Therefore, we now provide the results of our complementary security analysis, and estimate the success rate and average key rank metrics proposed in Paragraph A. As previously mentioned these evaluations are less confident, since for large  $q$  values such as  $q = 150$  we can have only one evaluation of the success function. Concretely, the best success rate / average key rank estimates are therefore obtained for  $q = 1$ . We took advantage of re-sampling when estimating them for larger  $q$ ’s.

Figures 6 and 7 illustrate these metrics are indeed correlated with the value of the PI estimates using the maximum profiling set, which explains the more efficient attacks against User 3. Concretely, the average rank figure suggests that correct PIN value can be exactly extracted in our 6-PIN case study with 5 to 10 observations for the most informative users and 30 to 40 observations for the least informative ones. The success rate curves also bring meaningful intuitions since they highlight that all (correct and random) PIN values can be correctly classified with our profiled models (in slightly more traces). This confirms our

<sup>4</sup> As mentioned in Section 2, this is due to the presence of another familiar event for this user, which he mentioned to us after the experiments were performed.

neurophysiological assumption from the previous section that the users react similarly to all random values.<sup>5</sup> Besides, Figure 6 is interesting since it shows how confidently the correct PIN value is classified independent of the others. Hence, its results would essentially scale with larger number of PIN values.

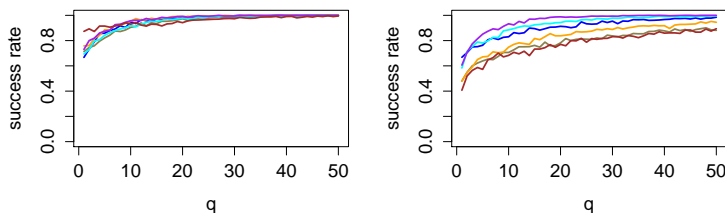


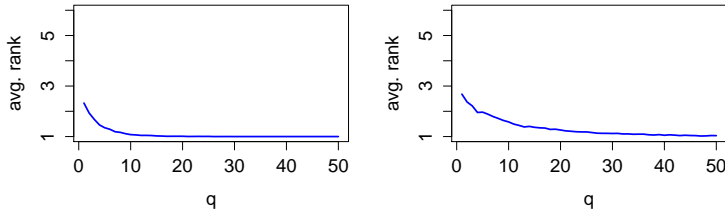
Fig. 6. Success rates per tag value for User 3 (left) and User 6 (right).

#### 4.2 Unsupervised (aka non-profiled) analysis

We now move to the more challenging problem of unsupervised / non-profiled attacks. For this purpose, we first applied the attack sketched in Section 3.3 with the maximum number of traces in the profiling set. That is, we repeated our evaluation of the PI metric six times, assuming each of the tag values to be the real one. Furthermore, we computed the confidence intervals for each of the PI estimates according to Section 3.2, Paragraph G. The results of this experiment are in Figure 8 for two users and lead to three observations.

First, looking at the first line of the figure, which corresponds to the correct PIN value, we can now confirm that the PI estimates of Section 4.1 are sufficiently accurate (e.g., the confidence intervals clearly guarantee a positive PI). Second, the confidence intervals for the random PIN values (i.e., tags 2 to 6) confirm the observation from our success rate curves (Figure 6) that the users react similarly to all random values. Third, the middle and bottom parts of the figure show the results of two (resp. 4) non-profiled attacks where the profiling set was split in 2 (resp. 4) independent parts (without re-sampling), therefore leading to the evaluation of 2 (resp. 4) confidence intervals for each tag value. As expected, it indicates that the information extraction is significantly more challenging in this unsupervised / non-profiled context. Concretely, the PI estimate for the correct PIN value consistently started to overlap with the ones of random PINs for all users, as soon as the number of attack traces  $q$  was below 200, and no clear gain for the correct PIN could be noticed below  $q = 100$ . This confirms the intuition that unsupervised / non-profiled side-channel attacks are generally more challenging than supervised / profiled ones (here, by an approximate factor 5 to 10 depending on the users).

<sup>5</sup> We may expect more singularities (such as the one of User 5) to appear and launch false alarms in case studies with more PIN values. Yet, this would not contradict the trend of a significantly reduced average rank for the correct PIN value.



**Fig. 7.** Avg. rank of the correct PIN for User 3 (left) and User 6 (right).

This conclusion also nicely matches the one in Section 4.1, Figure 5, where we already observed that the (offline) estimation of an informative model is more expensive than its (online) exploitation for PIN code recovery as measured by the success rate and average rank (by similar factors). Indeed, in the unsupervised / non-profiled context such an estimation has to be performed “on-the-fly”.

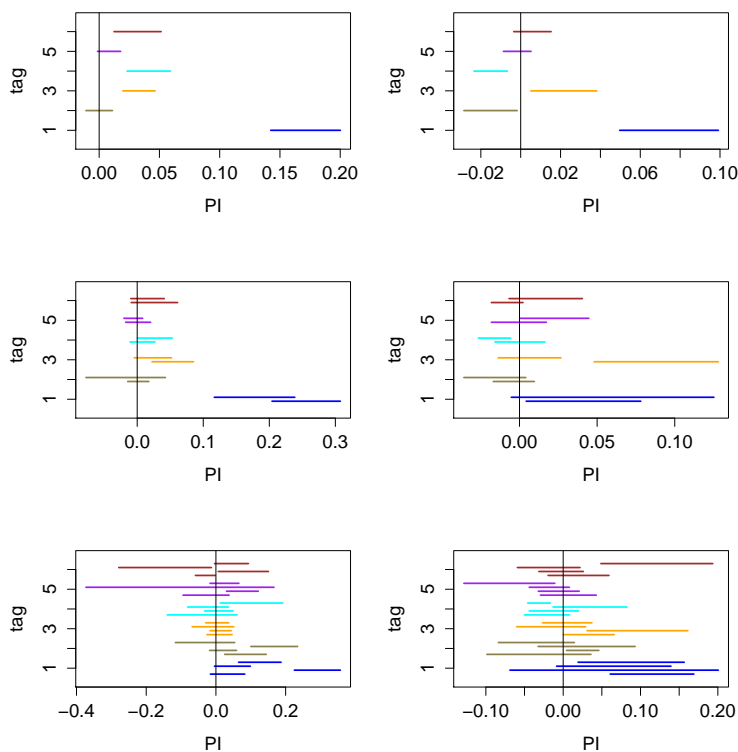
## 5 Consequences & Conclusions

The results in this paper lead to two important conclusions.

First, and from the security point-of-view, our experiments show that extracting concrete PIN codes from EEG signals, while theoretically feasible, may not be a very critical threat. PIN extraction attacks using BCIs indeed require several observations to succeed with high probability. Furthermore, the difference between the complexity of successful supervised / profiled attacks (around 10 correct PIN observations) and unsupervised / non-profiled attacks (more in the hundreds range) is noticeable. Yet, our results generally confirm the existence of exploitable information in EEG signals, which may become more worrying in case of targets with smaller cardinalities (e.g., extracting the knowledge of one relative among a set of unknown people displayed on a screen).

Second, and given the importance of profiling for efficient information extraction from EEG signals, our experiments underline that privacy issues may be even more worrying than security ones in BCI-based applications. Indeed, when it comes to privacy, the adversary trying to identify a user is less limited in his profiling abilities. In fact, any correlation between his target user and some feature found in a dataset is potentially exploitable. In this context, the data minimization principle does not seem to be a sufficient answer: it may be that the EEG signals collected for one (e.g., gaming) activity can be used to reveal various other types of (e.g., medical, political, ...) correlations. Anonymity is probably not the right answer either (since correlations with groups of users may be as discriminant as personal ones). And such issues are naturally amplified in case of malicious applications (e.g., it seem possible to design a BCI-based game where situations lead the users to incidentally reveal preferences). So overall, it appears as an important challenge to design tools that provide evidence of “fair treatment” when manipulating EEG signals, which can be connected to emerging challenges related to computations on encrypted data [25].





**Fig. 8.** Confidence intervals for the (non-profiled) PI evaluation of Section 3.3 with  $\approx 900$  observations (top),  $\approx 450$  observations (middle) and  $\approx 225$  observations (bottom), for Users 8 (left) and 6 (right).

## References

1. <http://emotiv.com/>, last retrieved july 2016.
2. <http://neurosky.com/>, last retrieved july 2016.
3. <http://www.chesworkshop.org/>, last retrieved july 2016.
4. C. Archambeau, E. Peeters, F. Standaert, and J. Quisquater. Template attacks in principal subspaces. In *CHES 2006. Proceedings*, volume 4249 of *LNCS*, pages 1–14. Springer, 2006.
5. L. Batina, B. Gierlichs, E. Prouff, M. Rivain, F. Standaert, and N. Veyrat-Charvillon. Mutual information analysis: a comprehensive study. *J. Cryptology*, 24(2):269–291, 2011.
6. I. Berlad and H. Pratt. P300 in response to the subject’s own name. *Electroencephalography and Clinical Neurophysiology*, 96(5):472–474, 1995.
7. T. Bonaci, R. Calo, and H. J. Chizeck. App stores for the brain : Privacy and security in brain-computer interfaces. *IEEE Technol. Soc. Mag.*, 34(2):32–39, 2015.
8. S. Chari, J. R. Rao, and P. Rohatgi. Template attacks. In *CHES. Proceedings.*, volume 2523 of *LNCS*, pages 13–28. Springer, 2002.

9. D. Coyle, J. C. Príncipe, F. Lotte, and A. Nijholt. Guest editorial: Brain/neuronal - computer game interfaces and interaction. *IEEE Trans. Comput. Intellig. and AI in Games*, 5(2):77–81, 2013.
10. A. Duc, S. Faust, and F. Standaert. Making masking security proofs concrete - or how to evaluate the security of any leaking device. In *EUROCRYPT 2015. Proceedings, Part I*, volume 9056 of *LNCS*, pages 401–429. Springer, 2015.
11. F. Durvaux, F. Standaert, and N. Veyrat-Charvillon. How to certify the leakage of a chip? In *EUROCRYPT. Proceedings*, volume 8441 of *LNCS*, pages 459–476. Springer, 2014.
12. B. Efron and R. J. Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
13. J. Engel, D. E. Kuhl, M. E. Phelps, and paul H. Crandall. Comparative localization of foci in partial epilepsy by pct and eeg. *Annals of Neurology*, 12(6):529–537, 1982.
14. L. A. Farwell and E. Donchin. The truth will out: Interrogative polygraphy (lie detection) with event-related brain potentials. *Psychophysiology*, 28(5):531–547, 1991.
15. M. Ienca. Hacking the brain: Brain–computer interfacing technology and the ethics of neurosecurity. *Ethics and Information Technology*, 18(2):117–129, 2016.
16. M. Inzlicht, I. McGregor, J. B. Hirsh, and K. Nash. Neural markers of religious conviction. *Psychological Science*, 20(3):385–392, 2009.
17. M. Kutas and S. A. Hillyard. Reading senseless sentences: Brain potentials reflect semantic incongruity. *Science*, 207:203–205, 1980.
18. M. Kutas and S. A. Hillyard. Brain potentials during reading reflect word expectancy and semantic association. *Nature*, 307:161–163, 1984.
19. C. Lin, R. Wu, S. Liang, W. Chao, Y. Chen, and T. Jung. Eeg-based drowsiness estimation for safety driving using independent component analysis. *IEEE Trans. on Circuits and Systems*, 52-I(12):2726–2738, 2005.
20. S. Mangard, E. Oswald, and T. Popp. *Power analysis attacks - revealing the secrets of smart cards*. Springer, 2007.
21. I. Martinovic, D. Davies, M. Frank, D. Perito, T. Ros, and D. Song. On the feasibility of side-channel attacks with brain-computer interfaces. In *USENIX Security Symposium. Proceedings*, pages 143–158. USENIX Association, 2012.
22. C. M. Portas, K. Krakow, P. Allen, O. Josephs, J. L. Armony, and C. D. Frith. Auditory processing across the sleep-wake cycle: simultaneous eeg and fmri monitoring in humans. *Neuron*, 28(3):991–999, 2000.
23. M. Renaud, F. Standaert, N. Veyrat-Charvillon, D. Kamel, and D. Flandre. A formal study of power variability issues and side-channel attacks for nanoscale devices. In *EUROCRYPT 2011. Proceedings*, volume 6632 of *LNCS*, pages 109–128. Springer, 2011.
24. B. W. Silverman. *Density Estimation for Statistics and Data Analysis*. CRC press, 1986.
25. N. P. Smart. Computing on encrypted data. *Kayaks & Dreadnoughts in a sea of crypto*, September 2016.
26. F. Standaert, F. Koeune, and W. Schindler. How to compare profiled side-channel attacks? In *ACNS. Proceedings*, volume 5536 of *LNCS*, pages 485–498, 2009.
27. F. Standaert, T. Malkin, and M. Yung. A unified framework for the analysis of side-channel key recovery attacks. In *EUROCRYPT. Proceedings*, volume 5479 of *LNCS*, pages 443–461. Springer, 2009.
28. N. Veyrat-Charvillon, B. Gérard, M. Renaud, and F. Standaert. An optimal key enumeration algorithm and its application to side-channel attacks. In *SAC. Proceedings*, volume 7707 of *LNCS*, pages 390–406. Springer, 2012.