

# Exploiting Ineffective Fault Inductions on Symmetric Cryptography

Christoph Dobraunig, Maria Eichlseder, Thomas Korak, Stefan Mangard,  
Florian Mendel and Robert Primas

Graz University of Technology, Austria

[first.last@iaik.tugraz.at](mailto:first.last@iaik.tugraz.at)

**Abstract.** Since the seminal work of Boneh et al., the threat of fault attacks has been widely known and new techniques for fault attacks and countermeasures have been studied extensively. The vast majority of the literature on fault attacks focuses on the ability of fault attacks to change an intermediate value to a faulty one, such as differential fault analysis (DFA), collision fault analysis, statistical fault attack (SFA), fault sensitivity analysis, or differential fault intensity analysis. The other aspect of faults—that faults can be induced and do not change a value—has been far less researched. In case of symmetric ciphers, this area is covered by ineffective fault attacks (IFA). However, IFA relies on the ability of an attacker to induce reproducible deterministic faults like stuck-at faults for a smaller intermediate structure (e.g., one bit or byte), which is often considered to be impracticable.

As a consequence, most countermeasures against fault attacks focus on the ability of faults to change intermediate values and usually try to detect such a change (detection-based), or to destroy the exploitable information if a fault happens (ineffective countermeasures). Such countermeasures implicitly assume that the release of “fault-free” ciphertexts in the presence of a fault-inducing attacker does not reveal any exploitable information. In this work, we challenge this assumption and show attacks that exploit the fact that intermediate values leading to such “fault-free” ciphertexts show a non-uniform distribution, while they should be uniformly distributed. The presented attacks are entirely practical and are demonstrated to work for software implementations of AES and for a hardware co-processor. These practical attacks rely on faults induced by means of clock glitches and hence, are achieved using only low-cost equipment. We target two countermeasures as example, simple time redundancy with comparison and an ineffective countermeasure presented at CHES 2014. However, our attacks can be applied to a wider range of countermeasures and are not restricted to these two countermeasures.

**Keywords:** fault attack · ineffective countermeasure · fault detection · countermeasure · statistical ineffective fault attack

## 1 Introduction

Shortly after the seminal work of Boneh et al. [9] showed fault attacks on RSA, it became clear that also symmetric schemes are susceptible to this type of active implementation attacks. Starting with the differential fault analysis (DFA) of DES by Biham and Shamir [6], a rich field of research emerged that focuses on techniques to recover the secret key from faulty ciphertexts. Meanwhile, there exists a wide range of publications on how to apply DFA attacks to different cryptographic algorithms and in particular to AES [2]. In addition, novel attack techniques have been introduced such as fault sensitivity analysis (FSA) [20], differential fault intensity analysis (DFIA) [16] and statistical fault attacks (SFA) [15].

In order to prevent an attacker from learning the secret key from faulty outputs, extensive research on countermeasures has been conducted. There are essentially two categories of countermeasures. The first category covers sensor-based countermeasures that aim at detecting the physical process of the fault induction. For example, protected implementations may include light, voltage and temperature sensors in order to detect fault inductions by lasers [26], voltage glitches [4], or temperature variations [19]. Countermeasures of this kind have a long tradition in smart card industry. However, as there are more and more ways to induce faults, the focus is more and more on countermeasures that aim at managing the effect of a fault induction. This is the second category of countermeasures and also the main focus of the academic research.

The effect of a fault induction on a cryptographic algorithm can be modelled as the change of an intermediate variable  $x$  to a faulty intermediate variable  $x'$ . Such a change of a value can occur due to a direct modification of the variable  $x$ , but also due to instruction skips or addressing errors. Independent of the exact effect that leads from  $x$  to  $x'$ , there are two approaches on how to prevent that this change is exploited by an attacker. The first approach is to detect the difference  $\Delta = x - x'$  by adding redundancy to a design and to suppress an output in case  $\Delta \neq 0$ . Corresponding redundancy techniques range from simple temporal or spacial duplications to error detection codes. The second approach for managing  $\Delta \neq 0$  are infection-based countermeasures. In this case, a cipher output is always provided, but the goal is to change the ciphertext in a way that the ciphertext becomes useless for an attacker.

While most attack techniques and countermeasures focus on exploiting or preventing information leakage in case  $\Delta \neq 0$ , the question of whether an attacker can also learn information from ineffective faults has not been explored in depth so far. A fault induction is ineffective in case the fault induction is performed (e.g., a voltage glitch is performed), but it holds that  $\Delta = 0$  and the cipher output is consequently not changed due to the fault induction. Information leakage on the secret key can occur in this case if there is a dependency between the fault induction being ineffective and the data that is processed in the device.

Exploiting ineffective fault inductions typically implies that a larger number of fault inductions needs to be performed than in case of classical fault attacks that exploit  $\Delta \neq 0$ . While in case of classical fault attacks a few selected fault inductions are usually sufficient to determine the key, the exploitation of ineffective faults requires that fault inductions are performed on many different data inputs in order to find ineffective fault inductions.

The first and, to the best of our knowledge, the only attacks that exploit ineffective inductions are ineffective fault attacks (IFA) by Clavier [10]. For IFA, it is assumed that an attacker can reliably force an intermediate value  $x$  to a known value (e.g., 0) by a fault induction. Then, the basic idea is to feed random input data into a device and to perform the fault inductions until an ineffective induction is observed. In this case, the intermediate  $x$  is known and the secret key can be determined. The big drawback of this attack is that it requires a precise fault induction for a large number of encryptions, which is very difficult to achieve in practice.

**Our contribution.** In this work, we generalize IFA attacks and introduce statistical ineffective fault attacks (SIFA). As we will argue and show with practical evaluations, SIFA is applicable whenever either SFA or IFA is applicable, but also in a broader range of scenarios beyond that, in particular in the presence of countermeasures. Our attack does not rely on a specific fault model. We simply require that there is some dependency between the observation of an ineffective fault induction and the faulted intermediate value  $x$ , though the attacker does not need to know any further details of this dependency. Concretely, this means that the probability for changing an intermediate value  $x$  due to a fault induction is not the same for all values  $x$ . This bias of the probabilities for ineffective

fault inductions is the sole requirement on the fault induction.

Like IFA, SIFA can be applied in settings where it is possible to perform many fault inductions on encryptions with different data inputs and to observe whether the fault induction was ineffective. While IFA typically requires strong fault models like stuck-at faults, the requirements for SIFA on the fault induction are minimal and corresponding faults can be induced easily in practice with a high frequency and without the need for sophisticated laboratory equipment.

To show this, we attack protected implementations that feature countermeasures against fault attacks like DFA and SFA. In particular, we target both a detection-based and an ineffective countermeasure. In fact, countermeasures that are based on managing a fault effect  $\Delta \neq 0$  are ideal targets for SIFA. These countermeasures allow the attacker to collect observations where the fault induction was ineffective. Our empirical study shows that these countermeasures can be easily bypassed in practice and that it is necessary to combine them with additional countermeasures in order to provide protection against SIFA attacks.

Our concrete attack results are as follows. First, we target a detection-based countermeasure for AES that uses simple time redundancy with subsequent comparison. In order to show the robustness of our attack this evaluation is performed on 3 different AES implementations, attacking 8-bit and 32-bit-bitsliced software implementations as well as a hardware co-processor. The fault is induced by using a simple clock glitch. In all cases, the number of needed faulty encryptions is comparably low. SFA is not applicable here since no exploitable faulty output is released. Although IFA is not prevented by simple time redundancy with subsequent comparison, it still relies on precise stuck-at faults in certain bytes, which are hard to achieve in practice, especially in the case of the 32-bit-bitsliced and the hardware co-processor implementations. In contrast, SIFA can exploit any case where ineffective faults lead to a biased distribution, even without knowledge about the distribution of these values.

We then target the ineffective countermeasure proposed at CHES 2014 by Tupsamudre et al. [27], where neither SFA nor IFA are applicable. Here, we extend the software AES implementation from the AVRCryptoLib [1] and evaluate our attack for multiple security parameterisations. Again, simple clock glitches are used to induce the required faults, resulting in attacks that are rather easy to execute in practice and do not require any expensive laboratory equipment.

**Related work.** SIFA extends and connects several other ideas that have previously been published in the literature. One key-point of the presented attack is the fact that it exclusively exploits cases where a fault does not change the result of the computation. Therefore, our attack shares a common reference point with safe-error attacks [28] and IFA [10]. In a safe-error attack, the value of an intermediate variable is changed (fault effect  $\Delta \neq 0$ ) and the knowledge whether the faulted value is used or not is exploited. Typically, safe-error attacks are used to attack asymmetric schemes. In contrast, ineffective fault attacks [10] exploit specific cases where  $\Delta = 0$  and the fault shows no effect. More concretely, IFA relies on strong and known fault models, like precise stuck-at-0 faults, in order to probe values of intermediate variables.

We extend this idea from stuck-at faults as already used by Biham and Shamir [6] to the case that ineffective faults lead to a non-uniform distribution of intermediate values. As a result, we do not probe specific values; rather, we probe distributions. Hence, we are naturally able to deal with noise (e.g., no faults induced, or faults induced at a wrong position), which allows us to demonstrate the attack in practice. The methods that we use to exploit non-uniform distributions are related to those in statistical fault attacks (SFA) [15].

**Outline.** First, we give a short overview and summary of statistical fault attacks and both reviewed countermeasures in section 2. Then, we state the idea and show the working principle of the attack in section 3. Section 4 contains the results of our practical attack and at last, we conclude the paper in section 5.

## 2 Background

In this section, we first give a brief introduction to countermeasures against fault attacks and review the two countermeasures we put our focus on more closely. Then we discuss two fault attacks which are related to our attack: ineffective fault attacks and statistical fault attacks.

### 2.1 Countermeasures

To protect against fault attacks, countermeasures aim to detect or prevent faults either on the physical layer (e.g., light sensors or supply voltage detectors) or on an algorithmic level. In this work, we solely focus on the second category. The strategy of detection-based countermeasures is to detect that a fault changes an intermediate value, e.g., by performing redundant operations. If a fault is detected, the computation is aborted and no ciphertext is returned. In contrast, infection-based countermeasures always return a ciphertext, but attempt to process the ciphertext in such a way that the output becomes useless for an attacker in case of faults during the computation. In the following we review the detection-based and the infection-based countermeasure that we target in our practical evaluation.

#### 2.1.1 Detection-based Countermeasure

In this work, we consider detection-based countermeasures that detect faults by means of redundant operations. An overview of various techniques that achieve detection of faults with the help of redundant operations is given by Bar-El et al. [4]. We focus on simple time redundancy with comparison, although the attack is applicable to a wide range of detection-based countermeasures. The idea of this countermeasure (Algorithm 1) is to encrypt each plaintext block twice. Then, the resulting ciphertexts are compared. Only if they match, the ciphertext is released.

---

**Algorithm 1** Simple time redundancy with comparison

---

**Input:** key  $K$ , plaintext  $P$

**Output:** ciphertext  $C = E_K(P)$ , or  $\perp$

---

```

1:  $C_1 \leftarrow E_K(P)$ 
2:  $C_2 \leftarrow E_K(P)$ 
3: if  $C_1 \neq C_2$  return  $\perp$ 
4: return  $C_1$ 

```

---

Detection-based countermeasures would also allow to include mechanisms that disable a device upon a certain amount of fault inductions. While this approach sounds very appealing at first glance, it is very hard to realize in practice. On the one hand, there is the risk of false positives that might lead to the disabling of a device in regular use cases (e.g., due to supply problems when being powered by an electromagnetic field). On the other hand, there is the need to count faults in such a way that it cannot be easily bypassed by an attacker. For example, a simple increasing of a counter in non-volatile memory can be detected easily by an attacker in the power trace and due to the timing

behavior. Hence, an attacker can detect whether a fault induction was effective or not and can remove the power supply during the programming of the memory in order to prevent the increasing of the counter. No sophisticated equipment is needed in this case. For more secure counting mechanisms, dedicated hardware support is required, which is not available in most devices (e.g., IoT devices) that are exposed to fault attacks. However, in applications that allow realizing fault counting in a secure and reliable manner without the risk of too many false positives, it is an effective countermeasures against classic fault attacks as well as IFA.

### 2.1.2 Infective Countermeasure

In contrast to detection-based countermeasures that aim to detect a fault and then do not release a ciphertext, infective countermeasures always provide a cipher output, but change the ciphertext in a way that the ciphertext becomes useless for an attacker. As an example for infection-based countermeasures, we consider the infective countermeasure presented by Tupsamudre et al. at CHES 2014 [27] as an extension of an infective countermeasure presented by Gierlichs et al. [17]. This countermeasure has been proven to be secure against differential fault analysis by Patranabis et al. [22] under the assumption that the sequence of executed instructions is neither skipped, nor altered. Thus, the only attacks on this countermeasure so far are attacks that either skip or alter instructions [5]. The approach is summarized in Algorithm 2.

---

**Algorithm 2** CHES 2014 infective countermeasure (algorithm taken from [27])

---

**Input:**  $P, k^j$  for  $j \in \{1, \dots, n\}$ ,  $(\beta, k^0)$ ,  $(n = 11)$  for AES-128  
**Output:**  $C = E_K(P)$ , or infected state

---

```

1: State  $R_0 \leftarrow P$ , Redundant state  $R_1 \leftarrow P$ , Dummy state  $R_2 \leftarrow \beta$ 
2:  $i \leftarrow 1, q \leftarrow 1$ 
3:  $rstr \leftarrow \{0, 1\}^t$  //  $\#1(rstr) = 2n, \#0(rstr) = t - 2n$ 
4: while  $q \leq t$  do
5:    $\lambda \leftarrow rstr[q]$  //  $\lambda = 0$  implies a dummy round
6:    $\kappa \leftarrow (i \wedge \lambda) \oplus 2(-\lambda)$ 
7:    $\zeta \leftarrow \lambda \cdot \lceil i/2 \rceil$  //  $\zeta$  is actual round counter, 0 for dummy
8:    $R_\kappa \leftarrow RoundFunction(R_\kappa, k^\zeta)$ 
9:    $\gamma \leftarrow \lambda(-i \wedge 1) \cdot BLFN(R_0 \oplus R_1)$  // check if  $i$  is even
10:   $\delta \leftarrow (-\lambda) \cdot BLFN(R_2 \oplus \beta)$ 
11:   $R_0 \leftarrow (-\gamma \vee \delta) \cdot R_0 \oplus ((\gamma \vee \delta) \cdot R_2)$ 
12:   $i \leftarrow i + \lambda$ 
13:   $q \leftarrow q + 1$ 
14: return  $R_0$ 

```

---

We will now give the basic intention behind Algorithm 2. For a more detailed description we refer to the original work of Tupsamudre et al. [27]. Algorithm 2 works on three different states  $R_0, R_1$  and  $R_2$ . State  $R_0$  is initialized with the plaintext  $P$  and is the state on which the primary AES computation is performed. State  $R_1$  is also initialized with  $P$  and serves as working state for the redundant AES computation. In the fault-free case, both states  $R_0$  and  $R_1$  should contain the ciphertext at the end of the computation. The state  $R_2$  is initialized with a random 128-bit value  $\beta$  and serves as working state for the dummy round calculations. The key  $k^0$  is chosen such that  $RoundFunction(\beta, k^0) = \beta$ .

Before the computation starts, a random string  $rstr$  of length  $t$  is initialized randomly so that it contains 22 bits “1” and  $t - 22$  bits “0”. The algorithm iterates over  $rstr$  and executes for every “1” an AES round on  $R_0$ , or a redundant round on  $R_1$  (22 rounds for 2 times 10 rounds AES plus 2 times the whitening key addition) in an alternating sequence,

i.e., if a round on  $R_0$  has been calculated, the next “1” executes a redundant round on  $R_1$  so that after this calculation, the content of  $R_0$  and  $R_1$  should be the same in a fault-free case. For every “0”, a dummy round is computed to update  $R_2$ . The security level with respect to the number of dummy rounds that are executed depends on the size of  $t$  and can be chosen by the developer.

After every executed AES round, the algorithm checks if any of the values in registers  $R_0$ ,  $R_1$ , or  $R_2$  has been modified ( $R_0 \neq R_1$  or  $R_2 \neq \beta$ ). If this is the case, state  $R_0$  is, from this point on, always overwritten with the content of  $R_2$ , which is then returned as ciphertext. Since the value stored in  $R_2$  is random and has never been mixed with, nor depends in any other way on the value of the secret key, learning this value should be useless for the attacker.

## 2.2 Statistical Fault Attacks

Statistical fault attacks (SFA) were introduced by Fuhr et al. [15] as a method to recover the secret key of AES if an attacker is able to change an intermediate variable to a biased (i.e., not uniformly distributed) value by inducing a fault. They considered three different fault models on byte level:

1. Stuck-at-0
2. Stuck-at-0 with probability 0.5, or logical AND with random uniform value with probability 0.5
3. Logical AND with random uniform value

Fuhr et al. [15] evaluated various key recovery strategies dependent on the round where the fault is induced. For instance, they showed that if the fault is induced in one byte right before the last MixColumns application, 6 faulty ciphertexts in case of fault model 1, 14 faulty ciphertexts in case of fault model 2, and 80 faulty ciphertexts in case of fault model 3 are needed to recover 4 bytes of the secret key. These attacks require to partially decrypt every ciphertext back to the faulted byte for each key candidate and measure the squared euclidean imbalance (SEI) of this byte. The key candidate that gives the highest SEI is most likely the correct one.

Since SFAs make use of the ability of faults to change an intermediate value (to a biased value), they can be prevented by both countermeasures discussed in subsection 2.1. To bypass both types of countermeasures, an attacker could try to induce identical faults in both redundant computations and thus evade detection. With strongly biased faults, this may be easier to achieve than for random faults, as has been demonstrated for a detection-based countermeasure [23]. However, the attacker’s task gets more and more complicated with increasing redundancy of the countermeasure.

## 2.3 Ineffective Fault Attacks

The idea of ineffective fault attacks (IFA) by Clavier [10] is that certain faults can be used to probe intermediate values of a cryptographic algorithm. This technique can be used to circumvent countermeasures like simple time redundancy with comparison. Consider an attacker who induces a stuck-at-0 fault in one byte during one execution of AES, while leaving the other one correct. If the attacker nevertheless receives an output (ciphertext), the faulted value must already have been 0 before the fault. If this stuck-at-0 fault is induced in one byte of the last AES round before the last key addition, we can immediately recover one byte of the last round key: All an attacker has to do is to guess one byte of the key, decrypt the corresponding byte of the correct ciphertext back to the intermediate byte that has been faulted, and check if the resulting byte value is 0. If it is 0, the guessed

key byte is the right one. This approach is applicable just as easily for more than two redundant computations, since only one computation needs to be faulted.

However, the assumption that an attacker is able to deterministically change the value of an intermediate variable to 0 requires a very strong and powerful attacker. In practice, an attacker is usually less powerful and has to consider, for instance, false positives in case of failed fault inductions that do not show any effect. For simple time redundancy with comparison, one solution for this specific problem would be to repeat the fault induction several times for encryptions of the same plaintext to get results which are more or less noise-free, as suggested by Clavier and Wurcker [11]. Using this strategy directly also causes troubles in the case of the CHES 2014 infective countermeasure [27]. Although the fault might be induced always at the same byte at the same time, an attacker does not know if the affected byte belongs to a dummy round or not.

In the following, we demonstrate that not only stuck-at faults can be exploited in IFA and introduce statistical ineffective fault attacks. On a high level, these attacks can be seen as an intersection of the principles exploited in the case of IFA [10] and SFA [15]. In section 3, we explain the necessary conditions for our attack to work and demonstrate in section 4 that they are usually fulfilled when attacking real devices with algorithmic countermeasures. In particular, the attacker does not need to assume any specific fault model and can successfully recover the key even with very “noisy” faults with unpredictable, unreliable effects.

### 3 Statistical Ineffective Fault Attack

In this section, we discuss the ideas behind the extension from ineffective fault attacks [10] (IFA) to statistical ineffective fault attacks (SIFA). First, we will review the effects of faults with the help of fault distribution tables to identify the necessary conditions for SIFA to work in subsection 3.1. Then we introduce the working principle of SIFA in subsection 3.2. Finally, we develop some theoretical background of our attacks in subsection 3.3.

#### 3.1 The Effects of Faults

The effects caused by faults during the execution of cryptographic primitives are manifold and depend on the method used to induce the fault (e.g., laser, clock glitches), the architecture and manufacturing technology of the attacked device, and various other parameters (e.g., targeting a register or arithmetic instruction). However, all faults have in common that they change the value of a  $b$ -bit intermediate variable from a value  $x$ , which it would have for the correct execution, to a value  $x'$  in the presence of a fault. Observing the probability of transitions from a certain value  $x \rightarrow x'$  gives us a fault distribution table (see subsection 3.3 for the exact definition).

With the help of such a fault distribution table, we are able to characterize the effects of a wide range of faults that can happen in practice. For example, this allows us to capture faults where the value of  $x'$  is independent of the value  $x$ , like stuck-at faults, random faults, and biased faults, but also more complex relations where  $x'$  depends in some sense on  $x$ , for instance by faulting the instruction that computes  $x$ . In Table 1, we show various examples of fault distribution tables for different faults on a 2-bit intermediate variable.

Most fault countermeasures that work on an algorithmic level can only conceal cases where  $x \neq x'$ , because a fault that results in  $x = x'$  is indistinguishable from a normal working condition. As a consequence, an attacker has access to ciphertexts where the attacked (faulted) intermediate variable follows a distribution determined by the diagonal (red values) in Table 1. The attacks presented in the following sections show that a non-uniform distribution in this diagonal can be exploited to recover the key. Therefore, for an implementation protected by such a fault countermeasure to be resistant against

Table 1: Fault distribution tables for several 2-bit fault models.

(a) Stuck-at-0					(b) Random-And						
		$x'$						$x'$			
		00	01	10	11			00	01	10	11
$x$	00	1	0	0	0	$x$	00	1	0	0	0
	01	1	0	0	0		01	$\frac{1}{2}$	$\frac{1}{2}$	0	0
	10	1	0	0	0		10	$\frac{1}{2}$	0	$\frac{1}{2}$	0
	11	1	0	0	0		11	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
(c) Bit-flip					(d) Random fault						
		$x'$						$x'$			
		00	01	10	11			00	01	10	11
$x$	00	0	0	0	1	$x$	00	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
	01	0	0	1	0		01	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
	10	0	1	0	0		10	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
	11	1	0	0	0		11	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

our attack, one of the two following conditions has to be fulfilled: Either the probability that an ineffective fault happens is negligible (as in Table 1c), or the distribution in the diagonal of the fault distribution table is uniform (as in Table 1d).

While in theory, the bit-flip and random fault models of Table 1c and Table 1d are not susceptible to SIFA, our practical experiments in section 4 indicate that countermeasures cannot rely on the hope that only such “perfect” fault models occur in practice. For instance, consider the case where a bit-flip occurs probabilistically with the tendency to flip more often from 1 to 0 than from 0 to 1, as illustrated in Table 2. The resulting distribution has a biased diagonal.

Table 2: Bit-flip from 1 to 0 with 75 % and from 0 to 1 with 50 %.

		$x'$			
		00	01	10	11
$x$	00	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
	01	$\frac{3}{8}$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{1}{8}$
	10	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{1}{8}$	$\frac{1}{8}$
	11	$\frac{9}{16}$	$\frac{3}{16}$	$\frac{3}{16}$	$\frac{1}{16}$

In the following section, we will explain how such distributions can be exploited. Since the fault distribution tables are typically not known by an attacker (unless the attacker is able to profile the device), our attack works without any knowledge of the fault distribution table. This is demonstrated by the practical attacks of section 4, which we perform without any knowledge of the underlying fault model and fault distribution table.

### 3.2 Working Principle

We now consider a protected AES implementation as an example to show the working principle of SIFA. The attack can be split into 3 phases. The first phase is the actual fault attack and collection of suitable ciphertexts. In the second phase, parts of the last round key are guessed and the distribution of an intermediate state is evaluated. In the last phase of the attack, the partial key-guesses are ranked according to a metric (e.g., the Squared Euclidean Imbalance (SEI)) and the correct key is identified.

**Collecting ciphertexts.** Assume that we target one byte before the last application of MixColumns. We request the ciphertexts for a number of plaintexts and fault each encryption. If the implementation is protected with a detection-based countermeasure, we only obtain those ciphertexts where the fault was ineffective; in case of an ineffective countermeasure, we need to filter for ineffective faults ourselves by comparing the obtained ciphertexts with a second, unfaulted encryption (or decryption).

**Key guessing.** Following the fault model of subsection 3.1, we obtain a set of filtered ciphertexts whose intermediate value in one byte before the last MixColumns is non-uniformly distributed according to the diagonal of the fault distribution table. This information can be exploited to recover 32 bits of the last round key with a key-guessing strategy similar to the approach of SFA [15]: The attacker guesses 4 bytes of the last round key  $K_{10}$  and partially decrypts the last operations for each correct ciphertext to obtain a partial state  $S_9$ :

$$S_9 = \text{MC}^{-1} \circ \text{SB}^{-1} \circ \text{SR}^{-1}(C \oplus K_{10}). \quad (1)$$

Then, the attacker can evaluate the distribution of the byte in (1) where the fault has been induced, for example by computing the Squared Euclidean Imbalance (SEI) of the byte for each key candidate. In case of evaluating the SEI, no information of the penultimate round key has to be guessed, because the constant key addition changes only the values of the byte, but has no influence on the non-uniformity of the distribution.

**Determining the correct key.** In the previous phase, for each key candidate, the SEI of the targeted byte has been calculated. We assume that the right key leads to the distribution with the highest SEI if a sufficient number of ciphertexts is evaluated. A closer insight in the number of needed ciphertexts is given in subsection 3.3.

We want to point out that this attack allows to exploit any ineffective fault that causes a non-uniform distribution of an intermediate value, even if the distribution is not known by an attacker (as demonstrated in section 4). In addition, SIFA is robust against noise introduced by failed fault induction attempts, or a fault induction in dummy rounds in the case of the ineffective countermeasure. In the next section, we provide a statistical model for our attacks and justify the use of the SEI.

### 3.3 Statistical Model

In this section, we provide a more detailed statistical model of the attack. Our aim is to investigate the effect of various parameters, such as the fault distribution and the configuration of the countermeasure, on the necessary number of faulted ciphertexts to perform the attack with a certain success probability. We compare two scenarios: The practical scenario where the fault distribution is unknown to the attacker (CHI/SEI statistic), but also the theoretical scenario where the attacker happens to know the distribution (LLR statistic). The emphasis of our analysis is on the hardest case: An unknown fault distribution, close to uniform, with additional noise induced by countermeasures.

We consider the  $b$ -bit intermediate variable which contains the result of the operation targeted by the fault, and consider its distribution during the attack in more detail. From the attacker’s point of view, the value of this variable on a particular input (in absence of faults) is a random variable  $X$  which depends on the input and key. Additionally, the random variable  $X'$  denotes the value of this variable on the same input, but where the attacker additionally attempted to fault the operation. We also refer to  $X$  and  $X'$  as “before” and “after” the fault, although this is not strictly accurate. Both  $X$  and  $X'$  take values  $x \in \mathcal{X} = \{0, \dots, 2^b - 1\}$ . The action of the fault can be characterized by the transition probabilities

$$p_x(x') := \mathbb{P}[X' = x' | X = x].$$

In practice, this fault distribution table FDT =  $(p_x(x'))_{x,x'}$  is usually not known, or can only be roughly estimated. To perform the proposed attack, the attacker does not need to know the FDT. However, the success and efficiency of the attack depends on some of the table’s properties. In the following, we will analyze the attack complexity and its dependency on the two relevant metrics: The fault’s ineffectivity rate  $\pi_{=}$ , and the capacity  $C(p)$  of the target distribution  $p$ .

### 3.3.1 Direct sampling: Detection countermeasure

We first consider attacks on detection-based countermeasures. We can only take advantage of samples where  $X = X'$ , i.e., the fault is ineffective. We assume that  $X$  is uniformly distributed, that is,  $\mathbb{P}[X = x] = 2^{-b}$ . Then, the probabilities  $\pi_{=}$  of an ineffective fault (ineffectivity rate) and  $\pi_{\neq}$  of an effective fault are

$$\pi_{=} = \mathbb{P}[X' = X] = \sum_{x' \in \mathcal{X}} \frac{p_{x'}(x')}{2^b}, \quad \pi_{\neq} = 1 - \pi_{=}.$$

We target the conditional distribution  $p_{=}(x')$  of  $X'$  in case of ineffective faults, i.e., the diagonal of the fault distribution table (see subsection 3.1):

$$p_{=}(x') := \mathbb{P}[X' = x' | X' = X] = \frac{p_{x'}(x')}{2^b \cdot \pi_{=}}.$$

The attacker neither knows this distribution, nor can she directly observe  $X'$ . However, based on the observed cipher output and a key hypothesis for the  $\kappa$ -bit last-round key material as in subsection 3.2, she obtains a hypothesis  $\hat{X}'$  for the value of  $X'$ , and can analyze the distribution  $\hat{p}$  of  $\hat{X}'$  for a fixed key guess across multiple samples. For an incorrect key guess, we assume a distribution very close to uniform<sup>1</sup>. For the correct key guess, we sample the unknown distribution  $p(x') = p_{=}(x')$ . If  $p_{=}(x')$  differs significantly from uniform, we can distinguish these two cases, and identify the samples from  $p_{=}(x')$  produced by the correct key  $k_0$  among the collection of samples from the nearly uniform distributions  $\theta_i(x') \approx \theta(x') = 2^{-b}$  produced by the wrong keys  $k_i$ ,  $1 \leq i < 2^\kappa$ .

To identify the correct key  $k_0$  and its distribution  $p = p_{=}$ , we associate a score statistic  $S(\hat{p})$  with each key candidate and the corresponding distribution  $\hat{p}$ , and rank the key candidates according to this statistic. This approach is closely related to statistical cryptanalysis, such as differential and linear cryptanalysis, and has been theoretically analyzed in those contexts. Under the assumption that  $S(\hat{p})$  is independently normally distributed for samples from either  $p$  or  $\theta$ ,

$$S(\hat{p}) \sim \begin{cases} \mathcal{N}(\mu_R, \sigma_R^2) & \text{if } \hat{p} \text{ was produced by } p, \\ \mathcal{N}(\mu_W, \sigma_W^2) & \text{if } \hat{p} \text{ was produced by } \theta, \end{cases} \quad (2)$$

<sup>1</sup>In practice, this is not necessarily the case, in particular for partially correct key guesses. For example, for a byte-stuck-at fault and a key guess that is only incorrect in one byte, the capacity is expected to drop from 255 to about 1, instead of 0.

Selçuk [25] analyzed the success probability of ranking the correct key  $k_0$  among the top  $2^{\kappa-a}$  of  $2^\kappa$  key candidates based on  $N$  samples, where  $a$  is the advantage. Then, the difference  $\Delta_a$  between the score of  $k_0$  and the score of the wrong key with rank  $2^{\kappa-a}$  (quantile  $\alpha = 1 - 2^{-a}$ ) is normally distributed with parameters

$$\Delta_a \sim \mathcal{N}(\mu_\Delta, \sigma_\Delta^2), \quad \begin{aligned} \mu_\Delta &= \mu_R - \mu_W - \sigma_W \Phi_{0,1}^{-1}(\alpha), \\ \sigma_\Delta^2 &\approx \sigma_R^2 \quad \text{for sufficiently large } 2^\kappa \text{ [8]}, \end{aligned}$$

and thus the success probability depending on  $N$  and  $a$  can be estimated as [25]

$$\mathbb{P}[\Delta_a > 0] \approx \Phi_{0,1} \left( \frac{\mu_R - \mu_W - \sigma_W \Phi_{0,1}^{-1}(\alpha)}{\sigma_R} \right). \quad (3)$$

To obtain useful complexity estimates from (3), we need a suitably distributed statistic  $S(\hat{p})$  and its parameters according to (2). We first consider the (unusual) case<sup>2</sup> that we know the real distribution  $p = p_-$ . Then, the Neyman-Pearson lemma [12, 21] states that the optimal statistic  $S$  is the log-likelihood ratio

$$S(\hat{p}) = \text{LLR}(\hat{p}) = \text{LLR}(\hat{p}, p, \theta) := N \sum_{x \in \mathcal{X}} \hat{p}(x) \log_2 \frac{p(x)}{\theta(x)}.$$

For large  $N$ ,  $\text{LLR}(\hat{p})$  tends towards a normal distribution as required in (2) [3, 12]. The success probability in (3) then depends on the Kullback-Leibler divergence  $D(p\|\theta)$ :

$$D(p\|\theta) := \sum_{\substack{x \in \mathcal{X} \\ p(x) \neq 0}} p(x) \log_2 \frac{p(x)}{\theta(x)}, \quad D_\Delta(p\|\theta) := \sum_{\substack{x \in \mathcal{X} \\ p(x) \neq 0}} p(x) \left[ \log_2 \frac{p(x)}{\theta(x)} \right]^2 - D(p\|\theta)^2.$$

If  $p$  is very close to uniform  $\theta$ , these can be approximated using the capacity  $C(p, \theta)$  [7]:

$$C(p, \theta) := \sum_{x \in \mathcal{X}} \frac{(p(x) - \theta(x))^2}{\theta(x)} \approx 2 D(p\|\theta) \approx D_\Delta(p\|\theta) \quad (\text{only if } p \text{ is close to } \theta.)$$

The resulting estimate for the necessary number of samples  $N_{\text{LLR}}$  to achieve a success probability  $P = \mathbb{P}[\Delta_a > 0]$  can be derived as [3, 8]:

$$N_{\text{LLR}} \approx \left[ \frac{\Phi_{0,1}^{-1}(P) \sqrt{D_\Delta(p\|\theta)} + \Phi_{0,1}^{-1}(\alpha) \sqrt{D_\Delta(\theta\|p)}}{D(p\|\theta) + D(\theta\|p)} \right]^2 \approx \frac{2[\Phi_{0,1}^{-1}(P) + \Phi_{0,1}^{-1}(\alpha)]^2}{C(p, \theta)}.$$

When applied to the AES fault analysis scenario, knowing  $p = p_-$  means both knowing the exact fault distribution of the ineffective faults and guessing the corresponding 8 key bits in the penultimate round, with a correspondingly increased advantage  $a$ . In the context of differential cryptanalysis, it has been demonstrated [8] that even small errors in the estimate of  $p$  can significantly increase the necessary number of samples. Since such exact models of  $p_-$  are usually not available for practical fault attacks, we can consider less optimal, but more robust statistics.

The classical test statistic for an unknown distribution  $p$  is Pearson's  $\chi^2$ :

$$S(\hat{p}) = \text{CHI}(\hat{p}) := \chi^2(\hat{p}, \theta) = N \sum_{x \in \mathcal{X}} \frac{(\hat{p}(x) - \theta(x))^2}{\theta(x)},$$

---

<sup>2</sup>Note that in this case, we could also target the last round with lower  $a$  and  $N$ , but more repetitions to obtain the full key.

or, for uniform  $\theta$ , the closely related Squared Euclidean Imbalance (SEI):

$$S(\hat{p}) = \text{SEI}(\hat{p}) := \sum_{x \in \mathcal{X}} (\hat{p}(x) - \theta(x))^2 = (N \cdot 2^b)^{-1} \cdot \text{CHI}(\hat{p}).$$

The statistic  $\text{CHI}(\hat{p})$  is distributed according to the (noncentral) chi-squared distribution with  $k = |\mathcal{X}| - 1 = 2^b - 1$  degrees of freedom and noncentrality parameter  $\lambda_{\text{R}} = N C(p, \theta)$  or  $\lambda_{\text{W}} = 0$ . For large  $k$  and  $N$ , this tends towards a normal distribution with parameters [14,18]

$$\text{CHI}(\hat{p}) \sim \begin{cases} \mathcal{N}(\mu_{\text{R}} = k + N C(p, \theta), & \sigma_{\text{R}}^2 = 2[k + 2N C(p, \theta)]), \\ \mathcal{N}(\mu_{\text{W}} = k, & \sigma_{\text{W}}^2 = 2k). \end{cases}$$

Based on these parameters, we can solve the quadratic equation in (3) to estimate the necessary number of samples as [8]

$$\begin{aligned} N_{\text{CHI}} &\approx \frac{s + \sqrt{s^2 - t}}{C(p, \theta)} & (s = \sqrt{2k} \Phi_{0,1}^{-1}(\alpha) + 2\Phi_{0,1}^{-2}(P), \quad t = 2k(\Phi_{0,1}^{-2}(\alpha) - \Phi_{0,1}^{-2}(P))) \\ &= \frac{\sqrt{2k} \Phi_{0,1}^{-1}(\alpha)}{C(p, \theta)} & (\text{for success probability } P = 0.5.) \end{aligned}$$

Summarizing, both statistics lead to an estimated number of samples that is proportional to  $1/C(p, \theta)$ , where the constant depends on the desired success probability  $P$  and advantage  $a$  (or quantile  $\alpha = 1 - 2^{-a}$ ). However, these estimates are only useful if the resulting  $N$  is reasonably large, that is, if  $p$  is not extremely different from  $\theta$ .

### 3.3.2 Noisy sampling: Infective countermeasure

So far, we assumed that for the correct key guess, the attacker makes the correct hypothesis  $\hat{X}' = X'$ , and thus directly samples the distribution  $p_{\pm}(x')$ . We will now show that the same approach also generalizes naturally to cases where the attacker only obtains noisy measurements.

As an example, consider the infective countermeasure with  $r$  dummy rounds. The attacker targets round  $R - t$  of the  $R = r + 11 + 11$  executed AES rounds, indexed  $1, \dots, R$ . To identify runs with ineffective faults, she has to compare the faulted ciphertexts  $C'$  with previously obtained correct ciphertexts  $C$  for the same plaintexts  $P$ , and keeps only the samples where  $C = C'$ . Assuming the same fault model as before, she will keep a fraction of about  $\pi_{\pm}$  samples. However, she does not know whether the ineffective fault really occurred in the penultimate AES round of the main (or, equivalently, redundant) encryption of  $P$ , or elsewhere: in a dummy round or the wrong AES round. The probability  $\sigma$  that the faulted round  $R - t$  was a relevant round depends on the attack setup. Figure 1 illustrates the practically observed probability ( $\sigma = \sigma_*$ , see subsection 4.1.3), as well as the expected probability if the attacker can equivalently target both main and redundant rounds in the practical attack setup ( $\sigma = \sigma_+$ ) or if she has to choose the target in advance ( $\sigma = \sigma_{\text{max}}$ ), where

$$\sigma = \begin{cases} \sigma_+ = \sigma_2 + \sigma_3, \\ \sigma_{\text{max}} = \max\{\sigma_2, \sigma_3\}, \end{cases} \quad \sigma_s = \frac{\binom{t}{s} \cdot \binom{R-t-1}{22-s-1}}{\binom{R}{22}}.$$

This probability is assumed to be independent of whether the fault was ineffective or not. Depending on whether round  $t$  was indeed relevant, the hypothesis  $\hat{X}'$  for the correct key now samples one of two distributions: If  $t$  was relevant,  $\hat{X}' = X'$  and we sample  $p_{\pm}(x')$ ; else, we sample a distribution close to uniform. Thus, we sample a noisy variable  $X''$  with distribution  $p_{\approx}(x'')$ , where

$$p_{\approx}(x'') = \sigma p_{\pm}(x'') + (1 - \sigma) 2^{-b} = \sigma (p_{\pm}(x'') - 2^{-b}) + 2^{-b}.$$

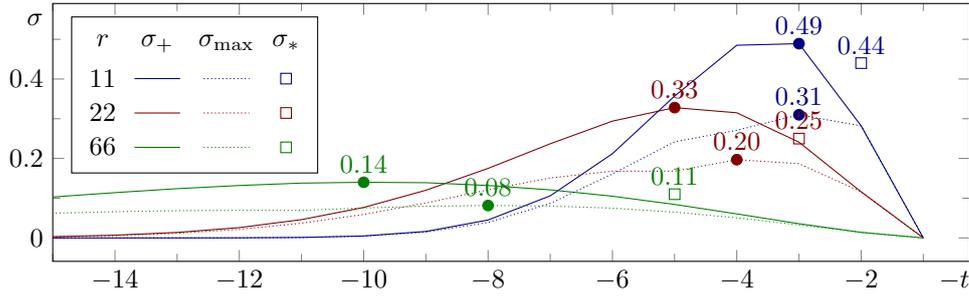


Figure 1: Probability  $\sigma$  of successful sampling for  $r \in \{11, 22, 66\}$  dummy rounds.

The capacity of this distribution is

$$C(p_{\approx}) = \sum_{x \in \mathcal{X}} \frac{(p_{\approx}(x) - 2^{-b})^2}{2^{-b}} = \sum_{x \in \mathcal{X}} \frac{(\sigma(p_{\approx}(x) - 2^{-b}))^2}{2^{-b}} = \sigma^2 C(p_{\approx}).$$

Thus, the expected data complexity for noisy sampling is  $\sigma^{-2}$  times higher compared to direct sampling.

In summary, the expected number of faults the attacker has to induce to collect enough samples is inverse proportional to  $\pi_{\approx} \cdot \sigma^2 \cdot C(p_{\approx})$ , where the constant depends on the desired success probability  $P$  and advantage  $a$ .

### 3.4 Examples and Simulations

To illustrate the statistical model in more detail, we consider a simulation of the attack with a random-and fault, i.e., each set bit of the target byte is flipped from 1 to 0 with probability  $\frac{1}{2}$ . The ineffectivity rate of this fault is  $\pi_{\approx} = (3/4)^8 \approx 10\%$ . We attack an AES implementation protected with the infective countermeasure (subsection 2.1) with  $r = 22$  dummy rounds and target round  $R - t = 44 - 4 = 40$ , obtaining a signal of  $\sigma = \frac{1111}{3526} \approx 0.315$  among the ineffectively faulted samples. The expected target distribution  $p(x)$  for the correct key is illustrated together with the uniform distribution  $\theta$  in Figure 2 and depends on the Hamming weight  $\text{hw}(x)$ :

$$p(x) = \sigma \cdot 2^{8-\text{hw}(x)} / 3^8 + (1 - \sigma) \cdot 2^{-8}.$$

To compare the practically necessary number of samples  $N$  (with or without knowledge of  $p(x)$ ) with the predictions of subsection 3.3, we evaluate the statistics  $\text{LLR}(\hat{p})$  and  $\text{CHI}(\hat{p})$  for the correct last-round key and for  $2^{24}$  wrong key candidates (out of  $2^{32}$ ; we set one byte to the correct value).

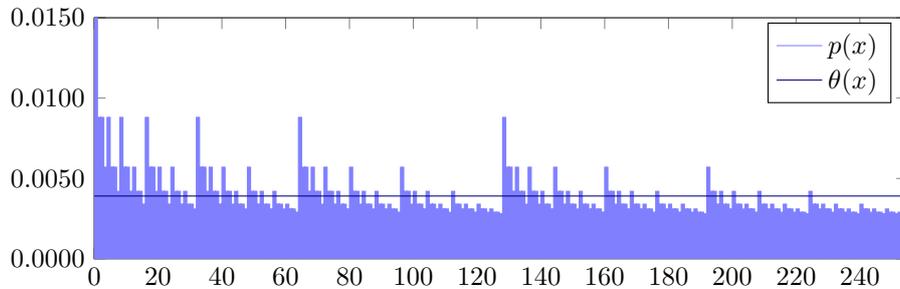


Figure 2: Distributions  $p$  and  $\theta$  for random-and fault and infective countermeasure.

For the  $\text{LLR}(\hat{p})$  statistic, we need to know the exact target distribution after addition of the penultimate round key, so we need to guess a byte  $K'$  of the penultimate round key in addition to the 24-bit key guess  $K$ . For simplicity, we evaluate each candidate  $K$  based on the statistic

$$S(\hat{p}) = \max_{K'} \text{LLR}(\hat{p}, p_{K'}, \theta) = \max_{K'} N \sum_{x \in \mathcal{X}} \hat{p}(x) \cdot \log_2 \frac{p(x \oplus K')}{\theta(x)}.$$

To reflect this in the model and evaluate the probability that the correct 24-bit  $K$  is ranked highest, we set the advantage to  $a = \kappa + 8 = 32$ , so  $\alpha = 1 - 2^{-32}$ . Based on the model of subsection 3.3, we expect the statistics  $\text{LLR}_R$  of the right key,  $\text{LLR}_W$  of any wrong key, and  $\text{LLR}_W^*$  of the best wrong key to be normally distributed with the following parameters:

$$\begin{aligned} \mu_R &= ND(p\|\theta) \approx 0.075N & \sigma_R^2 &= ND_{\Delta}(p\|\theta) \approx 0.252N \\ \mu_W &= -ND(\theta\|p) \approx -0.064N & \sigma_W^2 &= ND_{\Delta}(\theta\|p) \approx 0.157N \\ \mu_W^* &= \mu_W + \Phi_{0,1}^{-1}(\alpha) \sigma_W \approx -0.064N + 2.469\sqrt{N} & \sigma_W^{2*} &\ll \sigma_W^2. \end{aligned}$$

For the CHI statistic, we use  $a = \kappa = 24$  and expect the statistics  $\text{CHI}_R$ ,  $\text{CHI}_W$ , and  $\text{CHI}_W^*$  to be normally distributed with the following parameters:

$$\begin{aligned} \mu_R &= k + N C(p, \theta) \approx 255 + 0.131N & \sigma_R^2 &= 2k + 4N C(p, \theta) \approx 510 + 0.525N \\ \mu_W &= k = 255 & \sigma_W^2 &= 2k = 510 \\ \mu_W^* &= \mu_W + \Phi_{0,1}^{-1}(\alpha) \sigma_W \approx 375 & \sigma_W^{2*} &\ll \sigma_W^2. \end{aligned}$$

Figure 3 compares the resulting model (dashed:  $\mu_R, \mu_W^*$ ) with the statistics obtained in the practical key-recovery attack (solid:  $S(\hat{p}_R), S(\hat{p}_W^*)$ ). The predicted necessary number of samples  $N$  for success probability  $P = 0.8$  and with advantage  $a = 24$  (for CHI) or  $a = 32$  (for LLR) is marked as  $N_{\text{CHI}}$  and  $N_{\text{LLR}}$ , respectively. This estimate quite accurately matches the practically necessary  $N$ . It is worth noting that for both statistics, the best wrong key candidate scores slightly better than predicted with  $\mu_W$ . This can be partly explained with the not-entirely-uniform distribution of the target value for partially correct key guesses, as discussed in subsection 3.3. In case of  $\text{CHI}(\hat{p})$ , both the right and wrong keys scored slightly higher than expected.

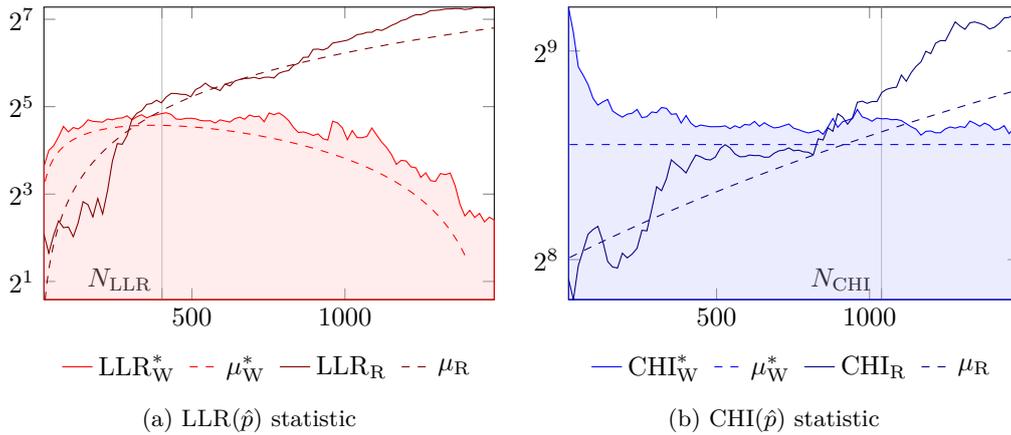


Figure 3: Simulation results for infection countermeasure and random-and fault.

We repeated simulations for other fault models (bit-stuck-at-0, byte-stuck-at-0) and for both countermeasures (detection, infective). The model matches the practical results

similarly well when the capacity  $C(p, \theta)$  is not too large and thus  $N$  is not too small to justify the normal approximation. In particular, for the byte-stuck-at-0 fault,  $C(p, \theta) \gg 1$ , and the model predicts fewer than the practically necessary  $N \approx 4$  ( $\sigma = 1$ , detection) or  $N \approx 15$  ( $\sigma = 0.315$ , infective) samples.

In summary, having insight in the concrete effect of a fault allows to model the scores of the correct and best wrong key and accurately predict the necessary number of samples for a successful attack. In such a case, the LLR outperforms the CHI (or SEI) statistic in terms of required samples. However, for LLR, even small errors in the estimate of  $p$  can significantly increase the necessary number of samples [8]. Hence, in practice, the CHI (or SEI) statistic is preferable, since the attacker can reliably and efficiently recover the key in the presence of countermeasures without any knowledge of  $p$  as demonstrated in section 4.

## 4 Practical Evaluation

For the practical evaluation of SIFA we have performed experiments targeting two types of fault countermeasures, implemented on various microcontrollers. The evaluated countermeasures are a detection-based and an infective countermeasure as described in Section 2.1. The target microcontrollers are listed in Table 3.

Table 3: Target microcontrollers of our attack evaluation

Name	ALU Size	Core	CPU Freq.
ATXmega 256A3	8-bit	Atmel AVR	12 MHz
ATXmega 128D4	8-bit	Atmel AVR	7 MHz
STM32 F3	32-bit	ARM Cortex-M4	7 MHz

In order to induce the faults we have used clock glitches. To be more precise, we xor an additional fast clock edge on the original clock signal to violate the critical path. By additionally varying width and offset of the induced clock edge, it is possible to influence the fault induction success rate and its impact on the faulted instruction. We have used an FPGA for generating both the original clock signal and the clock glitch for the device under test. For sake of simplicity, we determined our attack parameters with the help of an unprotected implementation in the case of the detection based countermeasure. In case that an unprotected implementation is not accessible to an attacker, determining the fault parameters is much more time consuming, but still feasible. We want to point out that the demonstrated attacks do not require a profiling of the actual distribution of the induced fault. In fact, the key recovery attacks have been performed without knowledge about the distribution of the targeted byte.

All experiments are performed in a fully automated attack setup. By using this setup, we are able to perform about 3 faulty encryptions per second, or 10 800 per hour. Depending on the concrete attack scenario, 1 000 to about 130 000 faulted encryptions are needed to reliably recover 4 bytes of the AES key. Hence, the time required to collect enough correct ciphertexts for key recovery is somewhere between 1 minute and 12 hours.

### 4.1 Attacks on Detection-based Countermeasure

At first we targeted a detection-based countermeasure that uses simple time redundancy with subsequent comparison (Algorithm 1). Here, the encryption is executed twice and only if the results of both encryptions are identical, the ciphertext is returned. Note that our attack is just as effective in case more than two redundant executions are performed. We evaluated our attack both for pure software AES implementation and AES

co-processor implementations. The attack against software AES was evaluated using an 8-bit register-based AES implementation on the ATXmega 256A3, and a 32-bit-bitsliced AES implementation on the STM32F3. The attack evaluation against the hardware co-processor AES implementation was performed on the ATXmega 256A3.

#### 4.1.1 8-bit Software AES on ATXmega 256A3

We used the ASM version of the AES-128 block cipher implementation from the AVRCryptoLib [1] as a basis for our protected AES implementation. Our attacks target the output of the S-box calculation of the 9<sup>th</sup> AES round. After experimenting with different fault parameterisations we have determined that inducing a clock glitch at 2044 clock cycles after the start of the encryption has the desired effect. We only induce a fault in one of the two AES encryptions and for approximately 34 % of the encryptions ciphertexts were received, indicating an ineffective fault. In total we performed 1000 encryptions, hence we received about 340 ciphertexts. The results presented in Figure 4 show that already after evaluating 230 correct ciphertexts, the SEI of the attacked byte is highest for the correct key candidate.

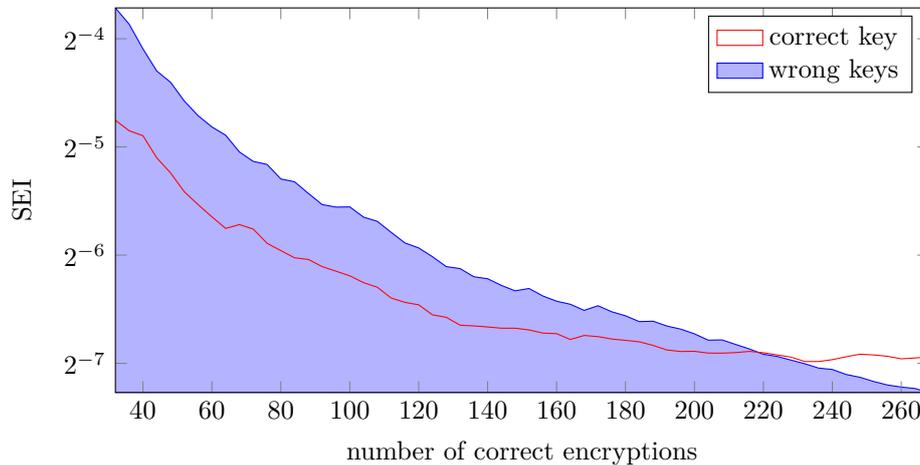


Figure 4: Attacks on software AES, ATXmega 256A3, detection countermeasure. SEI of the correct key ( $SEI_R$ ) vs. best SEI for a wrong key ( $SEI_W^*$ ) for  $N$  correct encryptions.

#### 4.1.2 32-bit-bitsliced Software AES on STM32F3

In order to evaluate our attack for bitsliced AES implementations, we have used the constant-time bitsliced implementation by Schwabe et al. [24]. The attack setup itself is similar to the one in the previous section. During the attack, we inserted a clock glitch 2463 cycles after the start of the encryption. In total we have performed 130000 faulted AES encryptions and received about 26000 correct ciphertexts. Hence, the fault probability is about 80 % in this setting. About 22000 correct ciphertexts were sufficient to reliably recover 4 bytes of the AES key, as seen in Figure 5.

#### 4.1.3 Hardware Co-Processor AES on ATXmega 256A3

For the fault attacks targeting the integrated AES co-processor on the ATXmega 256A3 microcontroller, we have used a clock glitch that is inserted 304 clock cycles after the start of the encryption. We induce a fault in one of the two AES encryptions, for approximately 69 % of the encryptions ciphertexts were received. For these encryptions,

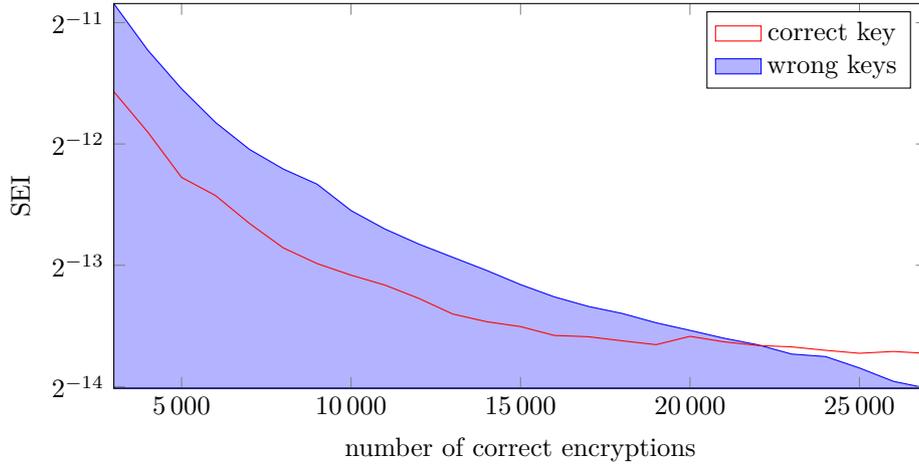


Figure 5: Attacks on bitsliced SW AES, STM32 F3, detection countermeasure. SEI of the correct key ( $SEI_R$ ) vs. best SEI for a wrong key ( $SEI_W^*$ ) for  $N$  correct encryptions.

the fault was ineffective. Similar to the evaluations of the software implementation, we performed 1 000 encryptions. In approximately 690 cases we received a ciphertext due to the fault probability of 31 % for this setting. The results presented in Figure 6 show that approximately 550 ciphertexts are sufficient for recovering 4 bytes of the AES key.

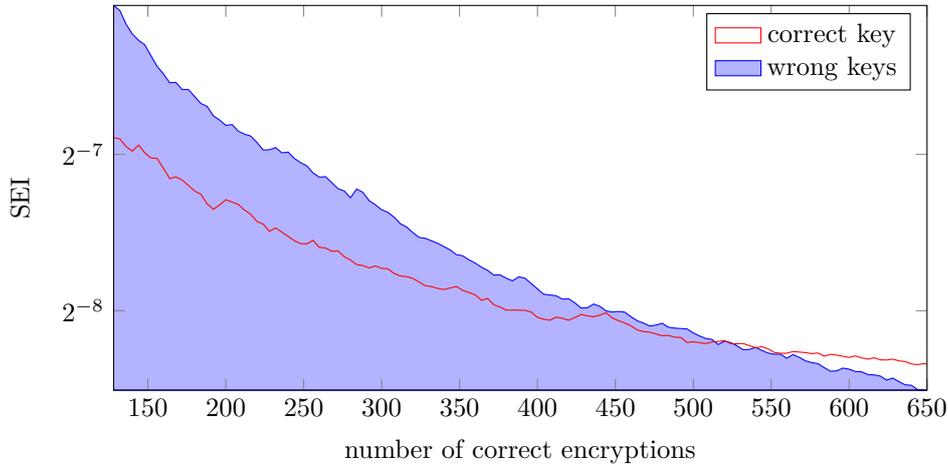


Figure 6: Attacks on HW AES co-processor, ATXmega 256A3, detection countermeasure. SEI of the correct key ( $SEI_R$ ) vs. best SEI of a wrong key ( $SEI_W^*$ ) for  $N$  correct encryptions.

## 4.2 Attacks on Infective Countermeasure

We then evaluated our attack on the infective countermeasure by Tupsamudre et al. [27] from CHES 2014 (Algorithm 2). Since the hardware co-processor of the ATXmega 256A3 only computes one complete call of AES, we limit this attack evaluation to purely software-based implementations.

First, we extended the C version of the AES-128 block cipher implementation from the AVRCryptoLib [1] according to Algorithm 2. The implementation of the AES round functions itself was not modified. The infective countermeasure requires 3 initial states

$R_0$ ,  $R_1$ , and  $R_2$ , two of which are set to the original plaintext,  $R_2$  is random. During encryption, 10 AES rounds are calculated for  $R_0$  and  $R_1$  in an alternating sequence. Additionally,  $t$  “dummy” rounds are calculated for  $R_2$  at randomly chosen points in time that are determined by *rstr*. Whenever a fault modifies a state during an AES round, the countermeasure amplifies the impact of the induced error up to a point where an exploitation of the faulty ciphertexts becomes infeasible. Thus, traditional differential fault attacks are not expected to work here. Still, by using the techniques introduced in Section 3.2, we show that such a countermeasure can be attacked by exploiting correct ciphertexts that are the result of biased but ineffective fault inductions.

Since the authors of the attacked countermeasure [27] did not give any recommendations for  $t$ , we have evaluated our attack for  $t = 11$ ,  $t = 22$ , and  $t = 66$ , leading to AES encryptions that require 33, 44, and 88 AES round function calls respectively. The resulting clock count per protected encryption is then about 110 000, 140 000, and 280 000.

The existence of dummy rounds makes the fault induction process less reliable. While we are able to induce faults in the detection-based scenario using constant trigger offsets with high success rates, the same cannot be said about the ineffective countermeasure scenario. Here, randomly occurring dummy rounds reduce the probability of hitting a penultimate AES round significantly.

We started with a simulation of multiple encryption runs in order to determine the round that performs the actual penultimate AES round with highest probability. Clearly, the best round for the attack depends on  $t$ . According to the simulation results in Table 4 we can expect the maximum possible fault induction success rate when targeting the 31<sup>st</sup>, 41<sup>st</sup>, and 83<sup>rd</sup> round. Once we know the best round for the attack, we can use a similar fault parameterisation as in the other experiments with the ATXmega 128D4. In contrast to the detection-based scenario, we cannot detect ineffective faults by observing just one encryption when infection is used. Hence, we always perform one encryption twice, once with fault induction, once without.

Table 4: Fault scenarios for attacking infective AES

Dummy Rounds ( $t$ )	Total AES Rounds	Target Round	Fault Offset (clks)	Maximum Possible Fault Induction Success
11	33	31	103 160	44%
22	44	41	129 440	25%
66	88	83	263 560	11%

The results of our attack evaluation are shown in Figure 7. Even though the low fault induction success rate results in a high number of correct ciphertexts, only very few of them will show a non-uniform distribution in the attacked byte. Thus the number of ciphertexts needed to distinguish a correct key candidate is significantly higher compared to the detection-based scenario. In the worst case, i.e.,  $t = 88$ , we had to perform about 130 000 faulty AES encryptions, which lead to about 120 000 correct ciphertexts. For a reliable recovery of 4 AES key bytes, about 110 000 such correct ciphertexts were necessary.

## 5 Conclusion

In this work, we provide an extensive insight on ineffective faults, where faults are being induced, but not showing an effect. The introduced statistical ineffective fault attacks (SIFA) can be seen as an intersection of the principles exploited by ineffective fault attacks (IFA) [10] and by statistical fault attacks (SFA) [15]. While previous work on IFA relies on strong models like stuck-at faults, we were able to relax these conditions up to a point were

we only require that intermediate values follow an unknown but non-uniform distribution. Hence, no special fault profiling of a targeted device is necessary.

SIFA inherits the ability from IFA that it only exploits the output of valid computations, which makes the attack independent of the degree of redundancy used in a countermeasure. As a consequence, it is not harder to attack a detection-based countermeasure performing 16 redundant operations compared to a countermeasure just performing 2. On the other hand, like SFA, SIFA works with minimal assumptions on the effect of the faults. Thus, similar as it has been shown for SFA (e.g., in [13]), we are able to demonstrate the feasibility of SIFA on various platforms in practice. However, in contrast to SFA, the practical attacks with SIFA are possible even in the presence of countermeasures against fault attacks.

We demonstrate the improvements of our work over IFA, amongst others by showing the applicability of SIFA on detection-based countermeasures utilizing 32-bit-bitsliced software AES implementations, or hardware co-processor AES implementations. In both cases the induction of precise stuck-at faults in certain bytes, as required by IFA, is considerably harder and was not possible in our fault setup.

Ultimately, we show that SIFA has new applications where neither SFA nor IFA are applicable, by demonstrating our attack on AES with the infective countermeasure that was presented at CHES 2014 [27]. Here, SFA does not work since the induced errors are amplified up to a point where the faulty output is unexploitable. Also IFA is not possible in this case, because even if precise faults were feasible, IFA cannot deal with the large amount of noise resulting from a low fault induction success rate that is caused by the presence of dummy rounds at random points in time. In addition, our experiments indicate that the number of dummy rounds  $t = 22$  seems to be a good choice. As shown in Figure 7,  $t = 22$  is significantly stronger than  $t = 11$ , while  $t = 66$  does not improve a lot compared to  $t = 22$ . Also in terms of performance,  $t = 22$  requires only 11 more AES rounds than  $t = 11$  while  $t = 66$  with 88 AES rounds in total appears to be an unreasonable additional cost, given the limited security gain.

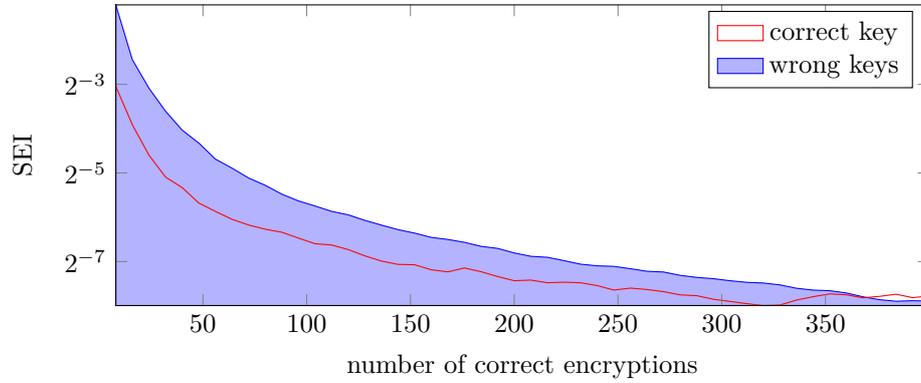
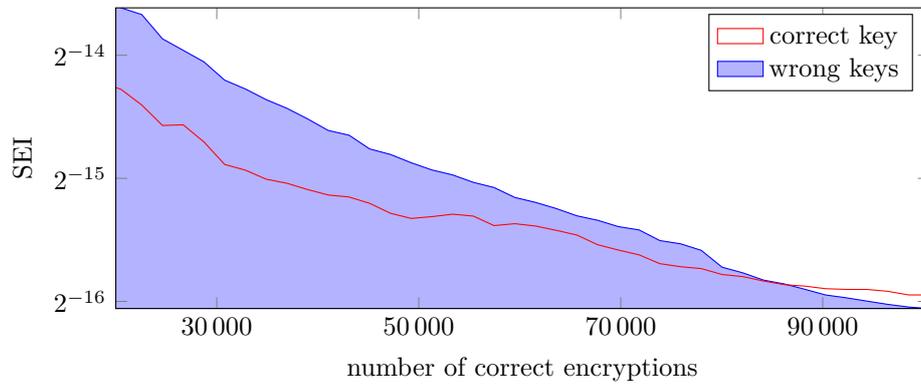
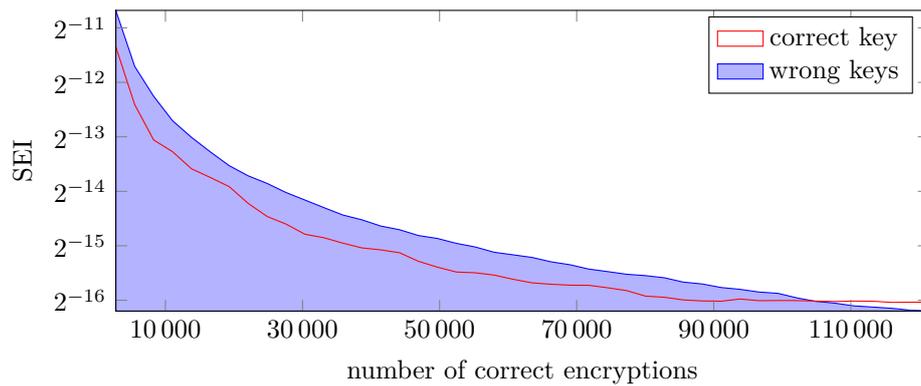
(a) Software AES, ATXmega 128D4,  $t = 33$ (b) Software AES, ATXmega 128D4,  $t = 44$ (c) Software AES, ATXmega 128D4,  $t = 88$ 

Figure 7: Attacks on software AES, ATXmega 128D4, infective countermeasure. SEI of the correct key ( $SEI_R$ ) vs. best SEI of a wrong key ( $SEI_W^*$ ) for  $N$  correct encryptions.

## References

- [1] AVR crypto lib. <https://git.cryptolib.org/?p=avr-crypto-lib.git>
- [2] Ali, S., Mukhopadhyay, D., Tunstall, M.: Differential fault analysis of AES: towards reaching its limits. *J. Cryptographic Engineering* 3(2), 73–97 (2013)
- [3] Baignères, T., Junod, P., Vaudenay, S.: How far can we go beyond linear cryptanalysis? In: Lee, P.J. (ed.) *Advances in Cryptology – ASIACRYPT 2004*. LNCS, vol. 3329, pp. 432–450. Springer (2004)
- [4] Bar-El, H., Choukri, H., Naccache, D., Tunstall, M., Whelan, C.: The sorcerer’s apprentice guide to fault attacks. *Proceedings of the IEEE* 94(2), 370–382 (2006)
- [5] Battistello, A., Giraud, C.: A note on the security of CHES 2014 symmetric infective countermeasure. In: Standaert, F.X., Oswald, E. (eds.) *Constructive Side-Channel Analysis and Secure Design – COSADE 2016*. LNCS, vol. 9689, pp. 144–159. Springer (2016)
- [6] Biham, E., Shamir, A.: Differential fault analysis of secret key cryptosystems. In: Kaliski Jr., B.S. (ed.) *Advances in Cryptology – CRYPTO ’97*. LNCS, vol. 1294, pp. 513–525. Springer (1997)
- [7] Biryukov, A., De Cannière, C., Quisquater, M.: On multiple linear approximations. In: Franklin, M.K. (ed.) *Advances in Cryptology – CRYPTO 2004*. LNCS, vol. 3152, pp. 1–22. Springer (2004)
- [8] Blondeau, C., Gérard, B., Nyberg, K.: Multiple differential cryptanalysis using LLR and  $\chi^2$  statistics. In: Visconti, I., Prisco, R.D. (eds.) *Security and Cryptography for Networks – SCN 2012*. LNCS, vol. 7485, pp. 343–360. Springer (2012)
- [9] Boneh, D., DeMillo, R.A., Lipton, R.J.: On the importance of checking cryptographic protocols for faults (extended abstract). In: Fumy, W. (ed.) *Advances in Cryptology – EUROCRYPT ’97*. LNCS, vol. 1233, pp. 37–51. Springer (1997)
- [10] Clavier, C.: Secret external encodings do not prevent transient fault analysis. In: Paillier, P., Verbauwhede, I. (eds.) *Cryptographic Hardware and Embedded Systems – CHES 2007*. LNCS, vol. 4727, pp. 181–194. Springer (2007)
- [11] Clavier, C., Wurcker, A.: Reverse engineering of a secret AES-like cipher by ineffective fault analysis. In: Fischer, W., Schmidt, J.M. (eds.) *Fault Diagnosis and Tolerance in Cryptography – FDTC 2013*. pp. 119–128. IEEE Computer Society (2013)
- [12] Cover, T.M., Thomas, J.A.: *Elements of Information Theory* (2. ed.). Wiley (2006)
- [13] Dobraunig, C., Eichlseder, M., Korak, T., Lomné, V., Mendel, F.: Statistical fault attacks on nonce-based authenticated encryption schemes. In: Cheon, J.H., Takagi, T. (eds.) *Advances in Cryptology – ASIACRYPT 2016*. LNCS, vol. 10031, pp. 369–395 (2016)
- [14] Drost, F.C., Kallenberg, W.C.M., Moore, D.S., Oosterhoff, J.: Power approximations to multinomial tests of fit. *Journal of the American Statistical Association* 84(405), 130–141 (1989)
- [15] Fuhr, T., Jaulmes, É., Lomné, V., Thillard, A.: Fault attacks on AES with faulty ciphertexts only. In: Fischer, W., Schmidt, J.M. (eds.) *Fault Diagnosis and Tolerance in Cryptography – FDTC 2013*. pp. 108–118. IEEE Computer Society (2013)

- [16] Ghalaty, N.F., Yuce, B., Taha, M.M.I., Schaumont, P.: Differential fault intensity analysis. In: *Fault Diagnosis and Tolerance in Cryptography – FDTC 2014*. pp. 49–58. IEEE Computer Society (2014)
- [17] Gierlichs, B., Schmidt, J.M., Tunstall, M.: Infective computation and dummy rounds: Fault protection for block ciphers without check-before-output. In: Hevia, A., Neven, G. (eds.) *Progress in Cryptology – LATINCRYPT 2012*. LNCS, vol. 7533, pp. 305–321. Springer (2012)
- [18] Hermelin, M., Cho, J.Y., Nyberg, K.: Multidimensional extension of Matsui’s Algorithm 2. In: Dunkelman, O. (ed.) *Fast Software Encryption – FSE 2009*. LNCS, vol. 5665, pp. 209–227. Springer (2009)
- [19] Hutter, M., Schmidt, J.: The temperature side channel and heating fault attacks. In: *Smart Card Research and Advanced Applications – CARDIS 2013*. LNCS, vol. 8419, pp. 219–235. Springer (2013)
- [20] Li, Y., Sakiyama, K., Gomisawa, S., Fukunaga, T., Takahashi, J., Ohta, K.: Fault sensitivity analysis. In: *Cryptographic Hardware and Embedded Systems – CHES 2010*. LNCS, vol. 6225, pp. 320–334. Springer (2010)
- [21] Neyman, P., Pearson, E.: On the problem of the most efficient tests of statistical hypotheses. *Philosophical Trans. of the Royal Society of London* pp. 289–337 (1933)
- [22] Patranabis, S., Chakraborty, A., Mukhopadhyay, D.: Fault tolerant infective countermeasure for AES. In: Chakraborty, R.S., Schwabe, P., Solworth, J.A. (eds.) *Security, Privacy, and Applied Cryptography Engineering – SPACE 2015*. LNCS, vol. 9354, pp. 190–209. Springer (2015)
- [23] Patranabis, S., Chakraborty, A., Nguyen, P.H., Mukhopadhyay, D.: A biased fault attack on the time redundancy countermeasure for AES. In: Mangard, S., Poschmann, A.Y. (eds.) *Constructive Side-Channel Analysis and Secure Design – COSADE 2015*. LNCS, vol. 9064, pp. 189–203. Springer (2015)
- [24] Schwabe, P., Stoffelen, K.: All the AES you need on cortex-m3 and M4. In: Avanzi, R., Heys, H.M. (eds.) *Selected Areas in Cryptography - SAC 2016*. LNCS, vol. 10532, pp. 180–194. Springer (2016)
- [25] Selçuk, A.A.: On probability of success in linear and differential cryptanalysis. *Journal of Cryptology* 21(1), 131–147 (2008)
- [26] Selmke, B., Brummer, S., Heyszl, J., Sigl, G.: Precise laser fault injections into 90 nm and 45 nm SRAM-cells. In: *Smart Card Research and Advanced Applications – CARDIS 2015*. LNCS, vol. 9514, pp. 193–205. Springer (2015)
- [27] Tupsamudre, H., Bisht, S., Mukhopadhyay, D.: Destroying fault invariant with randomization - A countermeasure for AES against differential fault attacks. In: Batina, L., Robshaw, M. (eds.) *Cryptographic Hardware and Embedded Systems – CHES 2014*. LNCS, vol. 8731, pp. 93–111. Springer (2014)
- [28] Yen, S.M., Joye, M.: Checking before output may not be enough against fault-based cryptanalysis. *IEEE Trans. Computers* 49(9), 967–970 (2000)