

# White-Box Cryptography: Don't Forget About Grey Box Attacks

Joppe W. Bos<sup>1</sup>, Charles Hubain<sup>2</sup>, Wil Michiels<sup>1,3</sup>, Cristofaro Mune<sup>4</sup>,  
Eloi Sanfelix Gonzalez<sup>4</sup>, and Philippe Teuwen<sup>2</sup>

<sup>1</sup> NXP Semiconductors

<sup>2</sup> Quarkslab

<sup>3</sup> Technische Universiteit Eindhoven

<sup>4</sup> Riscure

**Abstract.** Despite the fact that all current scientific white-box approaches of standardized cryptographic primitives have been publicly broken, these attacks require knowledge of the internal data representation used by the implementation. In practice, the level of implementation knowledge required is only attainable through significant reverse engineering efforts.

In this paper we describe new approaches to assess the security of white-box implementations which require *neither* knowledge about the look-up tables used *nor* expensive reverse engineering effort. We introduce the *differential computation analysis* (DCA) attack which is the software counterpart of the differential power analysis attack as applied by the cryptographic hardware community. Similarly, the *differential fault analysis* (DFA) attack is the software counterpart of fault-injection attacks on cryptographic hardware.

For DCA, we developed plugins to widely available dynamic binary instrumentation (DBI) frameworks to produce *software execution traces* which contain information about the memory addresses being accessed. For the DFA attack, we developed modified emulators and plugins for DBI frameworks that allow injecting faults at selected moments within the execution of the encryption or decryption process as well as a framework to automate static fault injection.

To illustrate the effectiveness, we show how DCA and DFA can extract the secret key from numerous publicly available non-commercial white-box implementations of standardized cryptographic algorithms. These approaches allow one to extract the secret key material from white-box implementations significantly faster and without specific knowledge of the white-box design in an automated or semi-automated manner.

## 1 Introduction

The widespread use of mobile “smart” devices enables users to access a large variety of ubiquitous services. This makes such platforms a valuable target (cf. [59])

---

This is an extended version of the article published by Springer-Verlag available at 10.1007/978-3-662-53140-2\_11 and includes the treatment of differential fault attacks.

for a survey on security for mobile devices). There are a number of techniques to protect the cryptographic keys residing on these mobile platforms. The solutions range from unprotected software implementations on the lower range of the security spectrum, to tamper-resistant hardware implementations on the other end. A popular approach which attempts to hide a cryptographic key inside a software program is known as a *white-box implementation*.

Traditionally, people used to work with a security model where implementations of cryptographic primitives are modeled as “black boxes”. In this black box model the internal design is trusted and only the in- and output are considered in a security evaluation. As pointed out by Kocher, Jaffe, and Jun [39] in the late 1990s, this assumption turned out to be false in many scenarios. This black-box may leak some meta-information: e.g., in terms of timing or power consumption. This side-channel analysis gave rise to the gray-box attack model. Since the usage of (and access to) cryptographic keys changed, so did this security model. In two seminal papers from 2002, Chow, Eisen, Johnson and van Oorschot introduce the white-box model and show implementation techniques which attempt to realize a white-box implementation of symmetric ciphers [22,21].

The idea behind the white-box attack model is that the adversary can be the owner of the device running the software implementation. Hence, it is assumed that the adversary has full control over the execution environment. This enables the adversary to, among other things, perform static analysis on the software, inspect and alter the memory used, and even alter intermediate results (similar to hardware fault injections). This white-box attack model, where the adversary is assumed to have such advanced abilities, is realistic on many mobile platforms which store private cryptographic keys of third-parties. White-box implementations can be used to protect which applications can be installed on a mobile device (from an application store). Other use-cases include the protection of digital assets (including media, software and devices) in the setting of digital rights management, the protection of payment credentials in Host Card Emulation (HCE) environments and the protection of credentials for authentication to the cloud. If one has access to a “perfect” white-box implementation of a cryptographic algorithm, then this implies one should not be able to deduce any information about the secret key material used by inspecting the internals of this implementation. This is equivalent to a setting where one has only black-box access to the implementation. As observed by [25] this means that such a white-box implementation should resist all existing and future side-channel attacks.

As stated in [21], “*when the attacker has internal information about a cryptographic implementation, choice of implementation is the sole remaining line of defense.*” This is exactly what is being pursued in a white-box implementation: the idea is to embed the secret key in the implementation of the cryptographic operations such that it becomes difficult for an attacker to extract information about this secret key even when the source code of the implementation is provided.

Note that this approach is different from anti-reverse-engineering mechanisms such as code obfuscation [8,43] and control-flow obfuscation [31] although

these are typically applied to white-box implementations as an additional line of defense.

On top of these, protocol level countermeasures can also be used to mitigate risks in cases in which the application is network-connected. In this work we focus exclusively on the robustness of *naked* white-box implementations, without any additional countermeasures.

Although it is conjectured that no long-term defense against attacks on white-box implementations exist [21], there are still a significant number of companies selling white-box solutions. It should be noted that there are almost no known published results on how to turn any of the standardized public-key algorithms into a white-box implementation, besides a patent by Zhou and Chow proposed in 2002 [74]. The other published white-box techniques exclusively focus on symmetric cryptography. However, all such published approaches have been theoretically broken (see Section 2 for an overview). A disadvantage of these published attacks is that they require detailed information on how the white-box implementation is constructed. For instance, knowledge about the exact location of the *S*-boxes or the round transitions might be required together with the format of the applied encodings to the look-up tables (see Section 2 on how white-box implementations are generally designed). Vendors of white-box implementations try to avoid such attacks by ignoring Kerckhoffs’s principle and keeping the details of their design secret (and change the design once it is broken).

**Our Contributions.** All current cryptanalytic approaches require detailed knowledge about the white-box design used: e.g. the location and order of the *S*-boxes applied and how and where the encodings are used. This preprocessing effort required for performing an attack is an important aspect of the value attributed to commercial white-box solutions. Vendors are aware that their solutions do not offer a long term defense, but compensate for this by, for instance, regular software updates. Our contributions are attacks that work in an automated way, and are therefore a major threat for the claimed security level of the offered solutions compared to the ones that are already known. For some of the attacks, we use dynamic binary analysis (DBA), a technique often used to improve and inspect the quality of software implementations, to access and control the intermediate state of the white-box implementation.

One approach to implement DBA is called dynamic binary instrumentation (DBI). The idea is that additional analysis code is added to the original code of the client program at run-time in order to aid memory debugging, memory leak detection, and profiling. The most advanced DBI tools, such as Valgrind [56] and Pin [44], allow one to monitor, modify and insert instructions in a binary executable. These tools have already demonstrated their potential for behavioral analysis of obfuscated code [63].

We have developed plugins for both Valgrind and Pin to obtain *software traces*<sup>1</sup>: traces which record the read and write accesses made to memory. Ad-

---

<sup>1</sup> The entire software toolchain ranging from the plugins, to the GUI, to the individual scripts to target the white-box challenges, to the tool to analyze the collected traces is released as open-source software: see <https://github.com/SideChannelMarvels>.

ditionally, we developed plugins and modified emulators to introduce faults into software, as well as a framework to automate static fault injection<sup>2</sup>. We introduce two new attack vectors that use these techniques in order to retrieve the secret key from a white-box implementation:

- *differential computation analysis* (DCA), which can be seen as the software counterpart of the differential power analysis (DPA) [39] techniques as applied by the cryptographic hardware community. There are, however, some important differences between the usage of the software and hardware traces as we outline in Section 3.2.
- *differential fault analysis* (DFA), which is equivalent to fault-injection [15,10] attacks as applied by the cryptographic hardware community, but uses software means in order to inject faults, which allows for dynamic but also static fault injection.

We demonstrate that DCA can be used to efficiently extract the secret key from white-box implementations which apply at most a single remotely handled external encoding. Similarly, we show that DFA can be applied to white-box implementations that do not apply external encoding to the output of the encryption or decryption process. We apply DFA and DCA techniques to the publicly available white-box challenges of standardized cryptographic algorithms we could find; concretely this means extracting the secret key from four white-box implementations of the symmetric cryptographic algorithms AES and DES.

In contrast to the current cryptanalytic methods to attack white-box implementations, these techniques do not require any knowledge about the implementation strategy used, can be mounted without much technical cryptographic knowledge in an automated way, and extract the key significantly faster. Besides this cryptanalytic framework we discuss techniques which could be used as countermeasures against DCA (see Section 3.4) and DFA (see Section 4).

The main reason why DCA works is related to the choice of (non-) linear encodings which are used inside the white-box implementation (cf. Section 2). These encodings do not sufficiently hide correlations when the correct key is used and enables one to run side-channel attacks (just as in gray-box attack model). Sasdrich, Moradi, and Güneysu looked into this in detail [60] and used the Walsh transform (a measure to investigate if a function is a balanced correlation immune function of a certain order) of both the linear and non-linear encodings applied in their white-box implementation of AES. Their results show extreme unbalance where the correct key is used and this explain why first-order attacks like DPA are successful in this scenario.

## 2 Overview of White-Box Cryptography Techniques

The white-box attack model allows the adversary to take full control over the cryptographic implementation and the execution environment. It is not surprising that, given such powerful capabilities of the adversary, the authors of the

---

<sup>2</sup> The static fault injection framework, the individual scripts and the tool to analyze the collected outputs are released in the same project.

original white-box paper [21] conjectured that no long-term defense against attacks on white-box implementations exists. This conjecture should be understood in the context of code-obfuscation, since hiding the cryptographic key inside an implementation is a form of code-obfuscation. It is known that obfuscation of *any* program is impossible [5], however, it is unknown if this result applies to a specific subset of white-box functionalities. Moreover, this should be understood in the light of recent developments where techniques using multilinear maps are used for obfuscation that may provide meaningful security guarantees (cf. [28,16,4]). In order to guard oneself in this security model in the medium- to long-run one has to use the advantages of a software-only solution. The idea is to use the concept of *software aging* [33]: this forces, at a regular interval, updates to the white-box implementation. It is hoped that when this interval is small enough, this gives insufficient computational time to the adversary to extract the secret key from the white-box implementation. This approach only makes sense if the sensitive data is only of short-term interest, e.g. the DRM-protected broadcast of a football match. However, the practical challenges of enforcing these updates on devices with irregular internet access should be noted.

Protocol level mitigations to limit the impact or applicability of these attacks could also be implemented. Additionally, risk mitigation techniques could be implemented to counter fraud in connected applications with a back-end component such as e.g. mobile payment solutions. However these mitigation techniques are outside the scope of this paper.

**External encodings.** Besides its primary goal to hide the key, white-box implementations can also be used to provide additional functionality, such as putting a fingerprint on a cryptographic key to enable traitor tracing or hardening software against tampering [49]. There are, however, other security concerns besides the extraction of the cryptographic secret key from the white-box implementation. If one is able to extract (or copy) the entire white-box implementation to another device then one has copied the functionality of this white-box implementation as well, since the secret key is embedded in this program. Such an attack is known as *code lifting*. A possible solution to this problem is to use external encodings [21]. When one assumes that the cryptographic functionality  $E_k$  is part of a larger ecosystem then one could implement  $E'_k = G \circ E_k \circ F^{-1}$  instead. The input ( $F$ ) and output ( $G$ ) encoding are randomly chosen bijections such that the extraction of  $E'_k$  does not allow the adversary to compute  $E_k$  directly. The ecosystem which makes use of  $E'_k$  must ensure that the input and output encodings are canceled. In practice, depending on the application, input or output encodings need to be performed locally by the program calling  $E'_k$ . E.g. in DRM applications, the server may take care of the input encoding remotely but the client needs to revert the output encoding to finalize the content decryption.

In this paper, we can mount successful attacks on implementations which apply *at most a single remotely handled external encoding*. When both the input is received with an external encoding applied to it remotely and the output is computed with another encoding applied to it (which is removed remotely) then the implementation is not a white-box implementation of a standard algorithm

(like AES or DES) but of a modified algorithm (like  $G \circ \text{AES} \circ F^{-1}$  or  $G \circ \text{DES} \circ F^{-1}$ ).

**General idea.** The general approach to implement a white-box program is presented in [21]. The idea is to use look-up tables rather than individual computational steps to implement an algorithm and to encode these look-up tables with random bijections. The usage of a fixed secret key is embedded in these tables. Due to this extensive usage of look-up tables, white-box implementations are typically orders of magnitude larger and slower than a regular (non-white-box) implementation of the same algorithm. It is common to write a program that automatically generates a random white-box implementation given the algorithm and the fixed secret key as input. The randomness resides in the randomly chosen bijections to hide the secret key usage in the various look-up tables.

In the remainder of this section we first briefly recall the basics of DES and AES, the two most commonly used choices for white-box implementations, before summarizing the scientific literature related to white-box techniques.

**Data Encryption Standard (DES).** The DES is a symmetric-key algorithm and published as a Federal Information Processing Standard (FIPS) for the United States in 1979 [70]. For the scope of this work it is sufficient to know that DES is an iterative cipher which consists of 16 identical rounds in a criss-crossing scheme known as a Feistel structure. One can implement DES by only working on 8-bit (a single byte) values and using mainly simple operations such as rotate, bitwise exclusive-or, and table lookups. Due to concerns of brute-force attacks on DES the usage of triple DES, which applies DES three times to each data block, has been added to a later version of the standard [70].

**Advanced Encryption Standard (AES).** In order to select a successor to DES, NIST initiated a public competition where people could submit new designs. After a roughly three year period the Rijndael cipher was chosen as AES [1,23] in 2000: an unclassified, publicly disclosed symmetric block cipher. The operations used in AES are, as in DES, relatively simple: bitwise exclusive-or, multiplications with elements from a finite field of  $2^8$  elements and table lookups. Rijndael was designed to be efficient on 8-bit platforms and it is therefore straight-forward to create a byte-oriented implementation. AES is available in three security levels. E.g. AES-128 is using a key-size of 128 bits and 10 rounds to compute the encryption of the input.

## 2.1 White-Box Results

**White-Box Data Encryption Standard (WB-DES).** The first publication attempting to construct a WB-DES implementation dates back from 2002 [22] in which an approach to create white-box implementations of Feistel ciphers is discussed. A first attack on this scheme, which enables one to unravel the obfuscation mechanism, took place in the same year and used fault injections [32] to extract the secret key by observing how the program fails under certain errors. In 2005, an improved WB-DES design, resisting this fault attack, was presented in [42]. However, in 2007, two differential cryptanalytic attacks [9] were presented

which can extract the secret key from this type of white-box [29,72]. This latter approach has a time complexity of only  $2^{14}$ .

**White-Box Advanced Encryption Standard (WB-AES).** The first approach to realize a WB-AES implementation was proposed in 2002 [21]. In 2004, the authors of [12] present how information about the encodings embedded in the look-up tables can be revealed when analyzing the lookup tables composition. This approach is known as the BGE attack and enables one to extract the key from this WB-AES with a  $2^{30}$  time complexity. A subsequent WB-AES design introduced perturbations in the cipher in an attempt to thwart the previous attack [18]. This approach was broken [55] using algebraic analysis with a  $2^{17}$  time complexity in 2010. Another WB-AES approach which resisted the previous attacks was presented in [73] in 2009 and got broken in 2012 with a work factor of  $2^{32}$  [54].

Another interesting approach is based on using the different algebraic structure for the same instance of an iterative block cipher (as proposed originally in [11]). This approach [34] uses dual ciphers to modify the state and key representations in each round as well as two of the four classical AES operations. This approach was shown to be equivalent to the first WB-AES implementation [21] in [40] in 2013. Moreover, the authors of [40] built upon a 2012 result [68] which improves the most time-consuming phase of the BGE attack. This reduces the cost of the BGE attack to a time complexity of  $2^{22}$ . An independent attack, of the same time complexity, is presented in [40] as well.

**Miscellaneous White-Box Results.** The above mentioned scientific work only relates to constructing and cryptanalyzing WB-DES and WB-AES. White-box techniques have been studied and used in a broader context. In 2007, the authors of [50] presented a white-box technique to make code tamper resistant. In 2008, the cryptanalytic results for WB-DES and WB-AES were generalized to any substitution linear-transformation (SLT) cipher [51]. In turn, this work was generalized even further and a general analytic toolbox is presented in [3] which can extract the secret for a general SLT cipher.

Formal security notions for symmetric white-box schemes are discussed and introduced in [61,25]. In [13] it is shown how one can use the ASASA construction with injective  $S$ -boxes (where ASA stands for the affine-substitution-affine [58] construction) to instantiate white-box cryptography. A tutorial related to white-box AES is given in [53].

## 2.2 Prerequisites of Existing Attacks

In order to put our results in perspective, it is good to keep in mind the exact requirements needed to apply the white-box attacks from the scientific literature. These approaches require at least a basic knowledge of the scheme which is white-boxed. More precisely, the adversary needs to

- know the type of encodings that are applied on *the intermediate results*,

- know which *cipher operations* are implemented by which (*network of*) *lookup tables*.

The problem with these requirements is that vendors of white-box implementations are typically reluctant in sharing any information on their white-box scheme (the so-called “security through obscurity”). If that information is not directly accessible but only a binary executable or library is at disposal, one has to invest a significant amount of time in reverse-engineering the binary manually. Removing several layers of obfuscation before retrieving the required level of knowledge about the implementations needed to mount this type of attack successfully can be cumbersome. This additional effort, which requires a high level of expertise and experience, is illustrated by the sophisticated methods used as described in the write-ups of the publicly available challenges as detailed in Section 3.3.

In contrast, the DCA and DFA approaches introduced in Sections 3 and 4 do not need to remove the obfuscation layers nor require significant reverse engineering of the binary executable.

### 3 Side Channel Analysis of White-Box Cryptographic implementations

#### 3.1 Differential Power Analysis

Since the late 1990s it is publicly known that the (statistical) analysis of a power trace obtained when executing a cryptographic primitive might correlate to, and hence reveal information about, the secret key material used [39]. Typically, one assumes access to the hardware implementation of a known cryptographic algorithm. With  $I(p_i, k)$  we denote a target intermediate state of the algorithm with input  $p_i$  and where only a small portion of the secret key is used in the computation, denoted by  $k$ . One assumes that the power consumption of the device at state  $I(p_i, k)$  is the sum of a data dependent component and some random noise, i.e.  $\mathcal{L}(I(p_i, k)) + \delta$ , where the function  $\mathcal{L}(s)$  returns the power consumption of the device during state  $s$ , and  $\delta$  denotes some leakage noise. It is common to assume (see e.g., [46]) that the noise is random, independent from the intermediate state and is normally distributed with zero mean. Since the adversary has access to the implementation he can obtain triples  $(t_i, p_i, c_i)$ . Here  $p_i$  is one plaintext input chosen arbitrarily by the adversary, the  $c_i$  is the ciphertext output computed by the implementation using a fixed unknown key, and the value  $t_i$  shows the power consumption over the time of the implementation to compute the output ciphertext  $c_i$ . The measured power consumption  $\mathcal{L}(I(p_i, k)) + \delta$  is just a small fraction of this entire power trace  $t_i$ .

The goal of an attacker is to recover the part of the key  $k$  by comparing the real power measurements  $t_i$  of the device with an estimation of the power consumption under all possible hypotheses for  $k$ . The idea behind a Differential Power Analysis (DPA) attack [39] (see [38] for an introduction to this topic) is

to divide the measurement traces in two distinct sets according to some property. For example, this property could be the value of one of the bits of the intermediate state  $I(p_i, k)$ . One assumes — and this is confirmed in practice by measurements on unprotected hardware — that the distribution of the power consumptions for these two sets is different (i.e., they have different means and standard deviations). In order to obtain information about part of the secret key  $k$ , for each trace  $t_i$  and input  $p_i$ , one enumerates all possible values for  $k$  (typically  $2^8 = 256$  when attacking a key-byte), computes the intermediate value  $g_i = I(p_i, k)$  for this key guess and divides the traces  $t_i$  into two sets according to this property measured at  $g_i$ . If the key guess  $k$  was correct then the difference of the subsets' averages will converge to the difference of the means of the distributions. However, if the key guess is wrong then the data in the sets can be seen as a random sampling of measurements and the difference of the means should converge to zero. This allows one to observe correct key guesses if enough traces are available. The number of traces required depends, among other things, on the measurement noise and means of the distributions (and hence is platform specific).

While having access to output ciphertexts is helpful to validate the recovered key, it is not strictly required. Inversely, one can attack an implementation where only the output ciphertexts are accessible, by targeting intermediate values in the last round. The same attacks apply obviously to the decryption operation.

The same technique can be applied on other traces which contain other types of side-channel information such as, for instance, the electromagnetic radiations of the device. Although we focus on DPA in this paper, it should be noted that there exist more advanced and powerful attacks. This includes, among others, higher order attacks [48], correlation power analyses [17] and template attacks [20].

### 3.2 Software Execution Traces

To assess the security of a binary executable implementing a cryptographic primitive, which is designed to be secure in the white-box attack model, one can execute the binary on a CPU of the corresponding architecture and observe its power consumption to mount a differential power analysis attack (see Section 3.1). However, in the white-box model, one can do much better as the model implies that we can observe everything without any measurement noise. In practice such level of observation can be achieved by instrumenting the binary or instrumenting an emulator being in charge of the execution of the binary. We chose the first approach by using some of the available Dynamic Binary Instrumentation (DBI) frameworks. In short, DBI usually considers the binary executable to analyze as the bytecode of a virtual machine using a technique known as just-in-time compilation. This recompilation of the machine code allows performing transformations on the code while preserving the original computational effects. These transformations are performed at the basic block<sup>3</sup> level and are stored

<sup>3</sup> A basic block is a portion of code with only one entry point and only one exit point. However, due to practical technicalities, the definition of a basic block Pin

in cache to speed up the execution. For example this mechanism is used by the Quick Emulator (QEMU, an open hypervisor that performs hardware virtualization) to execute machine code from one architecture on a different architecture, in this case the transformation is the architecture translation [7]. DBI frameworks, like Pin [44] and Valgrind [56], perform another kind of transformation: they allow to add custom callbacks in between the machine code instructions by writing plugins or tools which hook into the recompilation process. These callbacks can be used to monitor the execution of the program and track specific events. The main difference between Pin and Valgrind is that Valgrind uses an architecture independent Intermediate Representation (IR) called VEX which allows to write tools compatible with any architecture supported by the IR. We developed (and released) such plugins for both frameworks to trace execution of binary executables on x86, x86-64, ARM and ARM64 platforms and record the desired information: namely, the memory addresses being accessed (for read, write or execution) and their content. It is also possible to record the content of CPU registers but this would slow down acquisition and increase the size of traces significantly; we succeeded to extract the secret key from the white-box implementations without this additional information. This is not surprising as table-based white-box implementations are mostly made of memory look-ups and make almost no use of arithmetic instructions (see Section 2 for the design rationale behind many white-box implementations). In some more complex configurations e.g. where the actual white-box is buried into a larger executable it might be desired to change the initial behavior of the executable to call directly the block cipher function or to inject a chosen plaintext in an internal application programming interface (API). This is trivial to achieve with DBI, but for the implementations presented in Section 3.3, we simply did not need to resort to such methods.

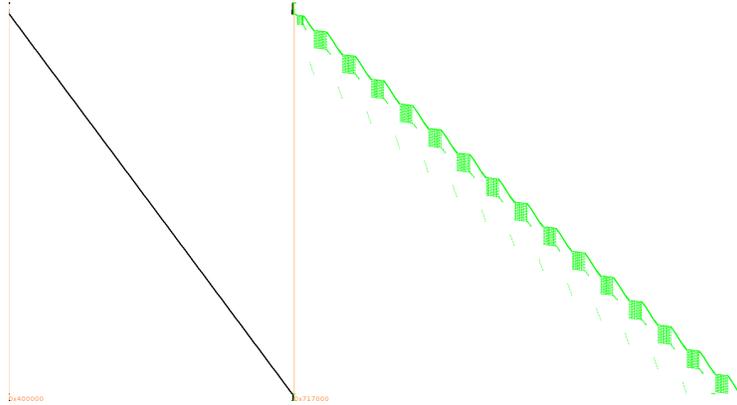
The following steps outline the process how to obtain software traces and mount a DPA attack on these software traces.

**First step.** Trace a single execution of the white-box binary with an arbitrary plaintext and record all accessed addresses and data over time. Although the tracer is able to follow execution everywhere, including external and system libraries, we reduce the scope to the main executable or to a companion library if the cryptographic operations happen to be handled there. A common computer security technique often deployed by default on modern operating systems is the Address Space Layout Randomization (ASLR) which randomly arranges the address space positions of the executable, its data, its heap, its stack and other elements such as libraries. In order to make acquisitions completely reproducible we simply disable the ASLR, as the white-box model puts us in control over the execution environment. In case ASLR cannot be disabled, it would just be a mere annoyance to realign the obtained traces.

**Second step.** Next, we visualize the trace to understand where the block cipher is being used and, by counting the number of repetitive patterns, determine

---

and Valgrind use is slightly different and may include several entry points or exit points.



**Fig. 1.** Visualization of a software execution trace of a white-box DES implementation.

which (standardized) cryptographic primitive is implemented: e.g., a 10-round AES-128, a 14-round AES-256, or a 16-round DES. To visualize a trace, we decided to represent it graphically similarly to the approach presented in [52]. Fig. 1 illustrates this approach: the virtual address space is represented on the  $x$ -axis, where typically, on many modern platforms, one encounters the text segment (containing the instructions), the data segment, the uninitialized data (BSS) segment, the heap, and finally the stack, respectively. The virtual address space is extremely sparse so we display only bands of memory where there is something to show. The  $y$ -axis is a temporal axis going from top to bottom. Black represents addresses of instructions being executed, green represents addresses of memory locations being read and red when being written. In Fig. 1 one deduces that the code (in black) has been unrolled in one huge basic block, a lot of memory is accessed in reads from different tables (in green) and the stack is comparatively so small that the read and write accesses (in green and red) are barely noticeable on the far right without zooming in.

**Third step.** Once we have determined which algorithm we target we keep the ASLR disabled and record multiple traces with random plaintexts, optionally using some criteria e.g. in which instructions address range to record activity. This is especially useful for large binaries doing other types of operations we are not interested in (e.g., when the white-box implementation is embedded in a larger framework). If the white-box operations themselves take a lot of time then we can limit the scope of the acquisition to recording the activity around just the first or last round, depending if we mount an attack from the input or output of the cipher. Focusing on the first or last round is typical in DPA-like attacks since it limits the portion of key being attacked to one single byte at once, as explained in Section 3.1. In the example given in Fig. 1, the read accesses pattern makes it trivial to identify the DES rounds and looking at the corresponding instructions (in black) helps defining a suitable instructions address range. While recording all memory-related information in the initial trace (first step), we only record a

single type of information (optionally for a limited address range) in this step. Typical examples include recordings of bytes being read from memory, or bytes written to the stack, or the least significant byte of memory addresses being accessed.

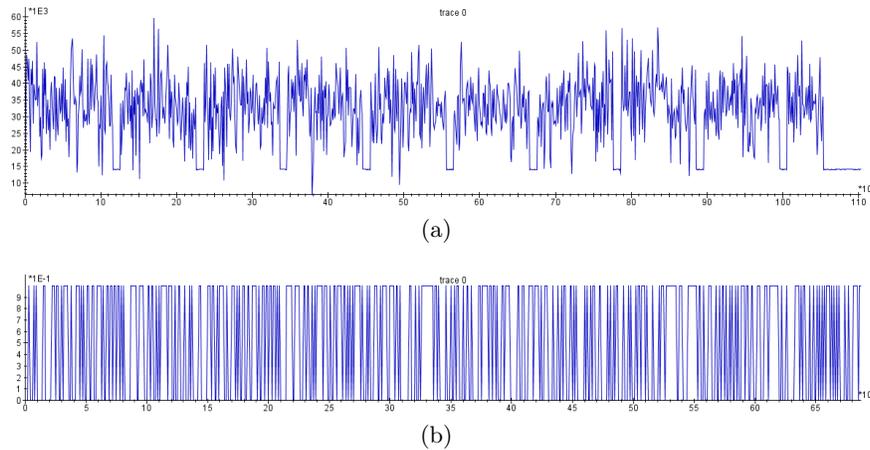
This generic approach gives us the best trade-off to mount the attack as fast as possible and minimize the storage of the software traces. If storage is not a concern, one can directly jump to the third step and record traces of the full execution, which is perfectly acceptable for executables without much overhead, as it will become apparent in several examples in Section 3.3. This naive approach can even lead to the creation of a fully automated acquisition and key recovery setup.

**Fourth step.** In step 3 we have obtained a set of software traces consisting of lists of (partial) addresses or actual data which have been recorded whenever an instruction was accessing them. To move to a representation suitable for usual DPA tools expecting power traces, we serialize those values (usually bytes) into vectors of ones and zeros. This step is essential to exploit all the information we have recorded. To understand it, we compare to a classical hardware DPA setup targeting the same type of information: memory transfers.

When using DPA, a typical hardware target is a CPU with one 8-bit bus to the memory and all eight lines of that bus will be switching between low and high voltage to transmit data. If a leakage can be observed in the variations of the power consumption, it will be an analog value proportional to the sum of bits equal to one in the byte being transferred on that memory bus. Therefore, in such scenarios, the most elementary leakage model is the Hamming weight of the bytes being transferred between CPU and memory. However, in our software setup, we know the exact 8-bit value and to exploit it at best, we want to attack each bit individually, and not their sum (as in the Hamming weight model). Therefore, the serialization step we perform (converting the observed values into vectors of ones and zeros) is as if in the hardware model each corresponding bus line was leaking individually one after the other.

When performing a DPA attack, a power trace typically consists of sampled analog measures. In our software setting we are working with *perfect* leakages (i.e., no measurement noise) of the individual bits that can take only two possible values: 0 or 1. Hence, our software tracing can be seen from a hardware perspective as if we were probing each individual line with a needle, something requiring heavy sample preparation such as chip decapping and Focused Ion Beam (FIB) milling and patching operations to dig through the metal layers in order to reach the bus lines without affecting the chip functionality. Something which is much more powerful and invasive than external side-channel acquisition.

When using software traces there is another important difference with traditional power traces along the time axis. In a physical side-channel trace, analog values are sampled at a fixed rate, often unrelated to the internal clock of the device under attack, and the time axis represents time linearly. With software execution traces we record information only when it is relevant, e.g. every time a byte is written on the stack if that is the property we are recording, and more-



**Fig. 2.** Figure (a) is a typical example of a (hardware) power trace of an unprotected AES-128 implementation (one can observe the ten rounds). Figure (b) is a typical example of a portion of a serialized software trace of stack writes in an AES-128 white-box, with only two possible values: zero or one.

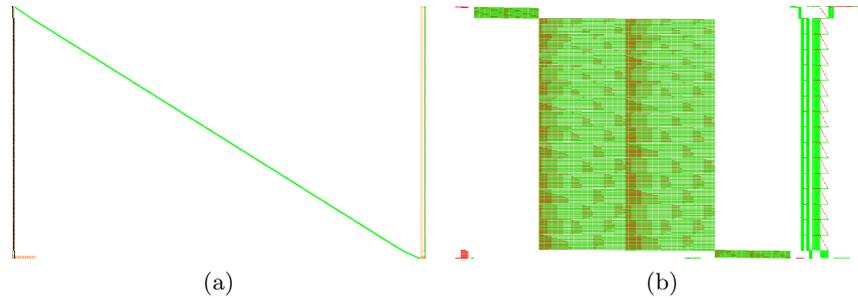
over bits are serialized as if they were written sequentially. One may observe that given this serialization and sampling on demand, our time axis does not represent an actual time scale. However, a DPA attack does not require a proper time axis. It only requires that when two traces are compared, corresponding events that occurred at the same point in the program execution are compared against each other. Figures 2a and 2b illustrate those differences between traces obtained for usage with DPA and DCA, respectively.

**Fifth step.** Once the software execution traces have been acquired and shaped, we can use regular DPA tools to extract the key. We show in the next section what the outcome of DPA tools look like, besides the recovery of the key.

**Optional step.** If required, one can identify the exact points in the execution where useful information leaks. With the help of *known-key correlation* analysis one can locate the exact “faulty” instruction and the corresponding source code line, if available. This can be useful as support for the white-box designer.

To conclude this section, here is a summary of the prerequisites of our differential computation analysis, in opposition to the previous white-box attacks’ prerequisites which were detailed in Section 2.2:

- Be able to run several times (a few dozens to a few thousands) the binary in a controlled environment.
- Having knowledge of the plaintexts (before their encoding, if any), or of the ciphertexts (after their decoding, if any).



**Fig. 3.** (a) Visualization of a software execution trace of the binary Wyseur white-box challenge showing the entire accessed address range. (b) A zoom on the stack address space from the software trace shown in (a). The 16 rounds of the DES algorithm are clearly visible.

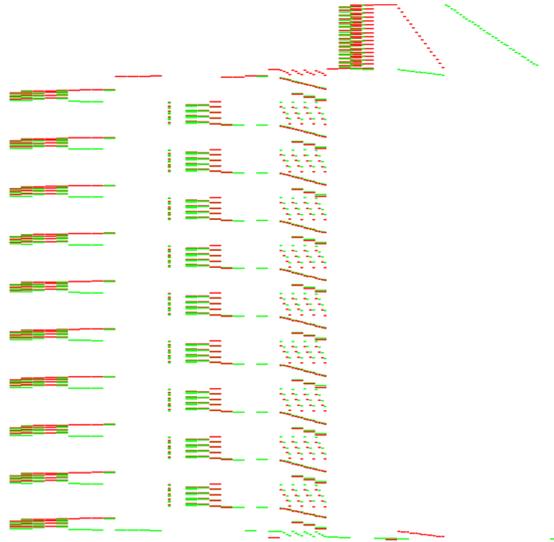
### 3.3 Analyzing Publicly Available White-Box Implementations

**The Wyseur Challenge** As far as we are aware, the first public white-box challenge was created by Brecht Wyseur in 2007. On his website<sup>4</sup> one can find a binary executable containing a white-box DES encryption operation with a fixed embedded secret key. According to the author, this WB-DES approach implements the ideas from [22,42] (see Section 2.1) plus “some personal improvements”. The interaction with the program is straight-forward: it takes a plaintext as input and returns a ciphertext as output to the console. The challenge was solved after five years (in 2012) independently by James Muir and “SysK”. The latter provided a detailed description [65] and used differential cryptanalysis (similar to [29,72]) to extract the embedded secret key.

Figure 3a shows a full software trace of an execution of this WB-DES challenge. On the left one can see the loading of the instructions (in black), since the instructions are loaded repeatedly from the same addresses this implies that loops are used which execute the same sequence of instructions over and over again. Different data is accessed fairly linearly but with some local disturbances as indicated by the large diagonal read access pattern (in green). Even to the trained eye, the trace displayed in Figure 3a does not immediately look familiar to DES. However, if one takes a closer look to the address space which represents the stack (on the far right) then the 16 rounds of DES can be clearly distinguished. This zoomed view is outlined in Figure 3b where the  $y$ -axis is unaltered (from Figure 3a) but the address-range (the  $x$ -axis) is rescaled to show only the read and write accesses to the stack.

Due to the loops in the program flow, we cannot just limit the tracer to a specific memory range of instructions and target a specific round. As a trace over the full execution takes a fraction of a second, we traced the entire program without applying any filter. The traces are easily exploited with DCA: e.g., if

<sup>4</sup> See <http://whiteboxcrypto.com/challenges.php>.



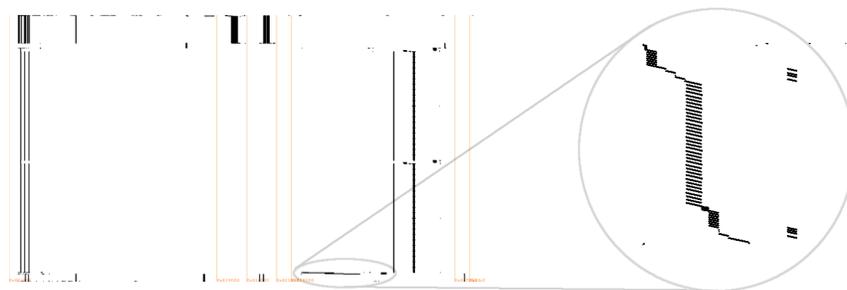
**Fig. 4.** Visualization of the stack reads and writes in a software execution trace of the Hack.lu 2009 challenge.

we trace the bytes written to the stack over the full execution and we compute a DPA over this entire trace without trying to limit the scope to the first round, the key is completely recovered with as few as 65 traces when using the output of the first round as intermediate value.

The execution of the entire attack, from the download of the binary challenge to full key recovery, including obtaining and analyzing the traces, took less than an hour as its simple textual interface makes it very easy to hook it to an attack framework. Extracting keys from different white-box implementations based on this design now only takes a matter of seconds when automating the entire process as outlined in Section 3.2.

**The Hack.lu 2009 Challenge** As part of the Hack.lu 2009 conference, which aims to bridge ethics and security in computer science, Jean-Baptiste Bédruone released a challenge [6] which consisted of a *crackme.exe* file: an executable for the Microsoft Windows platform. When launched, it opens a GUI prompting for an input, redirects it to a white-box and compares the output with an internal reference. It was solved independently by Eloi Vanderbéken [71], who reverted the functionality of the white-box implementation from encryption to decryption, and by “SysK” [65] who managed to extract the secret key from the implementation.

Our plugins for the DBI tools have not been ported to the Windows operating system and currently only run on GNU/Linux and Android. In order to use our tools directly we decided to trace the binary with our Valgrind variant and



**Fig. 5.** Visualization of the instructions in a software execution trace of the Karroumi WB-AES implementation by Klinec, with a zoom on the core of the white-box.

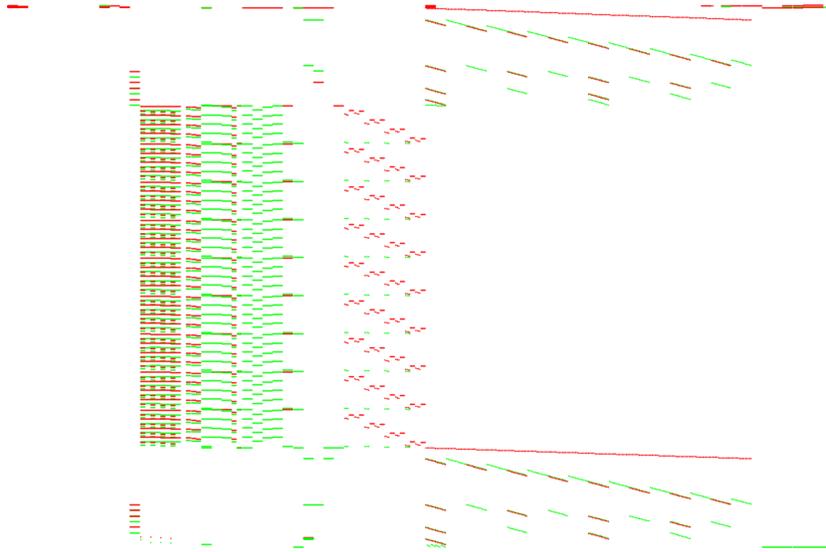
Wine [2], <sup>5</sup>, an open source compatibility layer to run Windows applications under GNU/Linux. We automated the GUI, keyboard and mouse interactions using `xdotool`<sup>6</sup>. Due to the configuration of this challenge we had full control on the input to the white-box. Hence, there was no need to record the output of the white-box and no binary reverse-engineering was required at all.

Fig. 4 shows the read and write accesses to the stack during a single execution of the binary. One can observe ten repetitive patterns on the left interleaved with nine others on the right. This indicates (with high probability) an AES encryption or decryption with a 128-bit key. The last round being shorter as it omits the *MixColumns* operation as per the AES specification. We captured a few dozen traces of the entire execution, without trying to limit ourselves to the first round. Due to the overhead caused by running the GUI inside Wine the acquisition ran slower than usual: obtaining a single trace took three seconds. Again, we applied our DCA technique on traces which recorded bytes written to the stack. The secret key could be completely recovered with only 16 traces when using the output of the first round *SubBytes* as intermediate value of an AES-128 encryption. As “SysK” pointed out in [65], this challenge was designed to be solvable in a couple of days and consequently did not implement any internal encoding, which means that the intermediate states can be observed directly. Therefore in our DCA the correlation between the internal states and the traced values get the highest possible value, which explains the low number of traces required to mount a successful attack.

**The SSTIC 2012 Challenge** Every year for the SSTIC, *Symposium sur la sécurité des technologies de l’information et des communications* (Information technology and communication security symposium), a challenge is published which consists of solving several steps like a Matryoshka doll. In 2012, one step of the challenge [47] was to validate a key with a Python bytecode “`check.pyc`”:

<sup>5</sup> <https://www.winehq.org/>

<sup>6</sup> <http://www.semicomplete.com/projects/xdotool/>



**Fig. 6.** Visualization of the stack reads and writes in the software execution trace portion limited to the core of the Karroumi WB-AES.

i.e. a marshalled object<sup>7</sup>. Internally this bytecode generates a random plaintext, forwards this to a white-box (created by Axel Tillequin) *and* to a regular DES encryption using the key provided by the user and then compares both ciphertexts. Five participants managed to find the correct secret key corresponding to this challenge and their write-ups are available at [47]. A number of solutions identified the implementation as a WB-DES without encodings (naked variant) as described in [22]. Some extracted the key following the approach from the literature while some performed their own algebraic attack.

Tracing the entire Python interpreter with our tool, based on either PIN or Valgrind, to obtain a software trace of the Python binary results in a significant overhead. Instead, we instrumented the Python environment directly. Actually, Python bytecode can be decompiled with little effort as shown by the write-up of Jean Sigwald. This contains a decompiled version of the “check.pyc” file where the white-box part is still left serialized as a pickled object<sup>8</sup>. The white-box makes use of a separate *Bits* class to handle its variables so we added some hooks to record all new instances of that particular class. This was sufficient. Again, as for the Hack.lu 2009 WB-AES challenge (see Section 3.3), 16 traces were enough to recover the key of this WB-DES when using the output of the first round as intermediate value. This approach works with such a low number of traces since the intermediate states are not encoded.

<sup>7</sup> <https://docs.python.org/2/library/marshal.html>

<sup>8</sup> <https://docs.python.org/2/library/pickle.html>

**Table 1.** DCA ranking for a Karroumi white-box implementation when targeting the output of the *SubBytes* step in the first round based on the least significant address byte on memory reads.

		key byte															
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
target bit	0	1	256	255	256	255	256	253	1	256	256	239	256	1	1	1	255
	1	1	256	256	256	1	255	256	1	1	5	1	256	1	1	1	1
	2	256	1	255	256	1	256	226	256	256	256	1	256	22	1	256	256
	3	256	255	251	1	1	1	254	1	1	256	256	253	254	256	255	256
	4	256	256	74	256	256	256	255	256	254	256	256	256	1	1	256	1
	5	1	1	1	1	1	1	50	256	253	1	251	256	253	1	256	256
	6	254	1	1	256	254	256	248	256	252	256	1	14	255	256	250	1
	7	1	256	1	1	252	256	253	256	256	255	256	1	251	1	254	1
All		✓	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓

**A White-Box Implementation of the Karroumi Approach** A white-box implementation of both the original AES approach [21] and the approach based on dual ciphers by Karroumi [34] is part of the Master thesis by Dušan Klinec [37]<sup>9</sup>. As explained in Section 2.1, this is the latest academic variant of [21]. Since there is no challenge available, we used Klinec’s implementation to create two challenges: one with and one without external encodings. This implementation is written in C++ with extensive use of the Boost<sup>10</sup> libraries to dynamically load and deserialize the white-box tables from a file. The left part of Figure 5 shows a software trace when running this white-box AES binary executable. The white-box code itself constitutes only a fraction of the total instructions; the right part of Figure 5 shows an enlarged view of the white-box core. Here, one can recognize the nine *MixColumns* operating on the four columns. This structure can be observed even better from the stack trace of Figure 6. Therefore we used instruction address filtering to focus on the white-box core and skip all the Boost C++ operations.

The best results were obtained when tracing the lowest byte of the memory addresses used in read accesses (excluding stack). Initially we followed the same approach as before: we targeted the output of the *SubBytes* in the first round. But, in contrast to the other challenges considered in this work, it was not enough to immediately recover the entire key. For some of the tracked bits of the intermediate value we observed a significant correlation peak: this is an indication that the first key candidate is very probably the correct one. Table 1 shows the ranking of the right key byte value amongst the guesses after 2000 traces, when sorted according to the difference of means (see Section 3.1). If the key byte is ranked at position 1 this means it was properly recovered by the attack. In total, for the first challenge we constructed, 15 out of 16 key bytes were ranked at position 1 for at least one of the target bits and one key byte

<sup>9</sup> The code be found at <https://github.com/ph4r05/Whitebox-crypto-AES>.

<sup>10</sup> <http://www.boost.org/>

**Table 2.** DCA ranking for a Karroumi white-box implementation when targeting the output of the multiplicative inversion inside the *SubBytes* step in the first round based on the least significant address byte on memory reads.

		key byte															
		0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
target bit	0	256	256	1	1	1	256	256	256	254	1	1	1	255	256	256	1
	1	1	1	253	1	1	256	249	256	256	256	226	1	254	256	256	256
	2	256	256	1	1	255	256	256	256	251	1	255	256	1	1	254	256
	3	254	1	69	1	1	1	1	1	252	256	1	256	1	256	256	256
	4	254	1	255	256	256	1	255	256	1	1	256	256	238	256	253	256
	5	254	256	250	1	241	256	255	3	1	1	256	256	231	256	208	254
	6	256	256	256	256	233	256	1	256	1	1	256	256	1	1	241	1
	7	63	256	1	256	1	255	231	256	255	1	255	256	255	1	1	1
All		✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

(key byte 6 in the table) did not show any strong candidate. However, recovering this single missing key-byte is trivial using brute-force.

Since the dual ciphers approach [34] uses affine self-equivalences of the original *S*-box, it might also be interesting to base the guesses on another target: the multiplicative inverse in the finite field of  $2^8$  elements (inside the *SubBytes* step) of the first round, before any affine transformation. This second attack shows results in Table 2 similar to the first one but distributed differently. With this sole attack the 16 bytes were successfully recovered — and the number of required traces can even be reduced to about 500 — but it may vary for other generations of the white-box as the distribution of leakages in those two attacks and amongst the target bits depends on the random source used in the white-box generator. However, when combining both attacks, we could always recover the full key.

It is interesting to observe in Table 1 and 2 that when a target bit of a given key byte does not leak (i.e. is not ranked first) it is very often *the worst* candidate (ranked at the 256<sup>th</sup> position) rather than being at a random position. This observation, that still holds for larger numbers of traces, can also be used to recover the key.

In order to give an idea of what can be achieved with an *automated* attack against new instantiations of this white-box implementation with other keys, we provide some figures: The acquisition of 500 traces takes about 200s on a regular laptop (dual-core i7-4600U CPU at 2.10GHz). This results in 832 kbits (104 kB) of traces when limited to the execution of the first round. Running both attacks as described in this section requires less than 30s. Attacking the second challenge with external encodings gave similar results. This was expected as there is no difference, from our adversary perspective, when applying external encodings or omitting them since in both cases we have knowledge of the original plaintexts before any encoding is applied.

**The NoSuchCon 2013 Challenge** In April 2013, a challenge designed by Eloi Vanderbéken was published for the occasion of the NoSuchCon 2013 conference<sup>11</sup>. The challenge consisted of a Windows binary embedding a white-box AES implementation. It was of “keygen-me” type, which means one has to provide a name and the corresponding *serial* to succeed. Internally the serial is encrypted by a white-box and compared to the MD5 hash of the provided name.

The challenge was completed by a number of participants (cf. [64,45]) but without ever recovering the key. It illustrates one more issue designers of white-box implementations have to deal with in practice: one can convert an encryption routine into a decryption routine without actually extracting the key.

For a change, the design is not derived from Chow [21]. However, the white-box was designed with external encodings which were *not* part of the binary. Hence, the user input was considered as encoded with an unknown scheme and the encoded output is directly compared to a reference. These conditions, without any knowledge of the relationship between the real AES plaintexts or ciphertexts and the effective inputs and outputs of the white-box, make it infeasible to apply a meaningful DPA attack, since, for a DPA attack, we need to construct the guesses for the intermediate values. Note that, as discussed in Section 2, this white-box implementation is *not* compliant with AES anymore but computes some variant  $E'_k = G \circ E_k \circ F^{-1}$ . Nevertheless we did manage to recover the key and the encodings from this white-box implementation with a new algebraic attack, as described in [67]. This was achieved after a painful de-obfuscation of the binary (almost completely performed by previous write-ups [64] and [45]), a step needed to fulfill the prerequisites for such attacks as described in Section 2.2.

The same white-box is found among the CHES 2015 challenges<sup>12</sup> in a Game-Boy ROM and the same algebraic attack is used successfully as explained in [66] once the tables got extracted.

### 3.4 Countermeasures against DCA

In hardware, counter-measures against DPA typically rely on a random source. The output can be used to mask intermediate results, to re-order instructions, or to add delays (see e.g. [19,30,62]). For white-box implementations, we cannot rely on a random source since in the white-box attack model such a source can simply be disabled or fixed to a constant value. Despite this lack of *dynamic* entropy, one can assume that the implementation which generates the white-box implementation has access to sufficient random data to incorporate in the generated source code and look-up tables. How to use this *static* random data embedded in the white-box implementation?

Adding (random) delays in an attempt to misalign traces is trivially defeated by using an address instruction trace beside the memory trace to realign traces automatically. In [24] it is proposed to use *variable* encodings when accessing the look-up tables based on the affine equivalences for bijective *S*-boxes (cf. [14] for

<sup>11</sup> See <http://www.nosuchcon.org/2013/>

<sup>12</sup> <https://ches15challenge.com/static/CHES15Challenge.zip>

algorithms to solve the affine equivalence problem for arbitrary permutations). As a potential countermeasure against DCA, the embedded (and possibly merged with other functionality) static random data is used to select which affine equivalence is used for the encoding when accessing a particular look-up table. This results in a variable encoding (at run-time) instead of using a fixed encoding. Such an approach can be seen as a form of masking as used to thwart classical first-order DPA.

One can also use some ideas from threshold implementations [57]. A threshold implementation is a masking scheme based on secret sharing and multi-party computation. One could also split the input in multiple shares such that not all shares belong to the same affine equivalence class. If this splitting of the shares and assignment to these (different) affine equivalence classes is done pseudo-randomly, where the randomness comes from the static embedded entropy and the input message, then this might offer some resistance against DCA-like attacks.

Another potential countermeasure against DCA is the use of external encodings. This was the primary reason why we were not able to extract the secret key from the challenge described in Section 3.3. However, typically the adversary can obtain knowledge related to the external encoding applied when he observes the behavior of the white-box implementation in the entire software-framework where it is used (especially when the adversary has control over the input parameters used or can observe the final decoded output). We stress that more research is needed to verify the strength of such approaches and improve the ideas presented here.

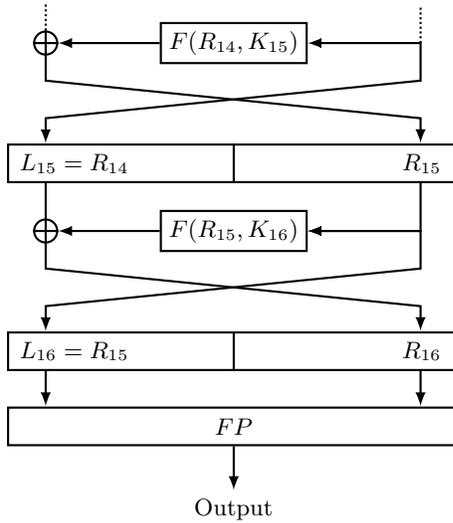
In practice, one might resort to methods to make the job of the adversary more difficult. Typical software counter-measures include obfuscation, anti-debug and integrity checks. It should be noted, however, that in order to mount a successful DCA attack one does not need to reverse engineer the binary executable. The DBI frameworks are very good at coping with those techniques and even if there are efforts to specifically detect DBI [27,41], DBI becomes stealthier too [36].

## 4 Differential Fault Analysis on White-box Cryptographic implementations

### 4.1 Differential Fault Analysis

Differential Fault Analysis was first introduced in [10], which presented the attack we describe in this section. The mechanics behind this type of attacks are as follows.

**First step.** The attacker repeats the target cryptographic algorithm with the same input multiple times and records correct and faulty outputs. In order to generate faulty outputs, the attacker introduces faults during the computation of the cryptographic algorithm, typically towards the end of its execution (e.g. before the final round). In the case of hardware solutions, this is typically done



**Fig. 7.** The final rounds of the DES cipher.

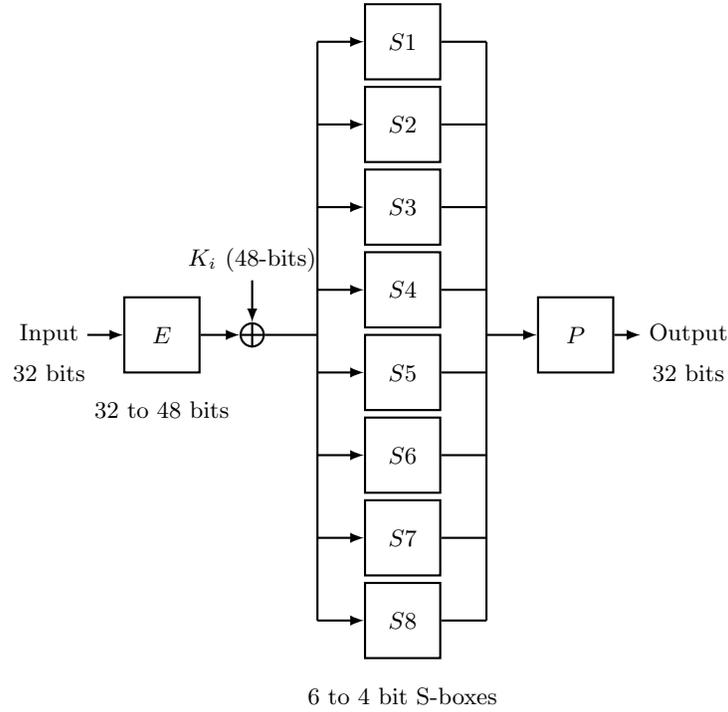
by altering the environmental conditions of the device. For example, voltage glitching or laser fault injection could be used to introduce faults.

**Second step.** The attacker builds a model of the introduced faults, and performs an analysis of the collected outputs (correct and faulty) in order to determine the secret key used.

The location within the cryptographic algorithm in which faults are injected in step 1 and the analysis performed in step 2 above are closely related. These depend on the cryptographic algorithm under attack. Attacks for a number of algorithms can be found in the literature, including but not limited to DES, AES, RSA and several ECC algorithms. We recall the differential fault attacks on triple DES and AES and show how to apply them to white-box implementations.

**DES DFA** Fig. 7 describes the final rounds of execution of a triple DES encryption or decryption, where the F function is defined as shown in Fig. 8. Each round of execution uses a 48-bit round key. The attack works by recovering one round key at a time, until the complete key can be computed.

In order to recover the last round key ( $K_{16}$ ) using a DFA attack, the attacker can inject faults during the execution of round 15 i.e. the computation of  $R_{15}$ .



**Fig. 8.** The Feistel function ( $F$ ) within the DES cipher.

Using the correct and faulty outputs, we have the following

$$\begin{aligned}
 L_{16} &= R_{15} \\
 R_{16} &= F(R_{15}, K_{16}) \oplus L_{15} \\
 \text{Output} &= FP(L_{16} \| R_{16}) \\
 L'_{16} &= R'_{15} \\
 R'_{16} &= F(R'_{15}, K_{16}) \oplus L_{15} \\
 \text{Output}' &= FP(L'_{16} \| R'_{16})
 \end{aligned}$$

In these expressions, the values of  $K_{16}$  and  $L_{15}$  are unknown. Combining these equations, we obtain the following equation in which the only unknown is the round key  $K_{16}$ :

$$R_{16} \oplus R'_{16} = F(R_{15}, K_{16}) \oplus F(R'_{15}, K_{16})$$

From the  $F$  function one can see that DES uses 6 round key bits to compute 4 output bits in each round, making it possible to solve this equation in chunks of 6 bits of the round key. In particular, for each individual S-Box the following equation needs to be solved:

$$(P^{-1}(R_{16} \oplus R'_{16}))_i = S_i(E(R_{15}) \oplus K_{16_i}) \oplus S_i(E(R'_{15}) \oplus K_{16_i})$$

Where  $E$  and  $P$  represent the expansion and permutation steps of the  $F$  function, respectively. This equation can be solved by exhaustive search.

Typically this results in a number of candidates for each affected sub-key per fault. Therefore, an attacker needs to iterate this process until only one candidate remains for each sub-key. However, in some cases, when the faults are not injected in the exact way as expected by the attack it may happen that a correct key gets discarded. Therefore, a counting strategy is used instead of discarding key candidates: for each fault, we compute the set of solutions to the equation above and increase the count for the respective candidates. Once all faults are analyzed, the candidate with the highest count for each sub-key is selected as the correct candidate. Once the last round key is known, the attack can be iterated to the previous round key. For this, the attacker injects faults one round earlier and computes the output of the one but last round by using the recovered last round key. If a triple DES cipher is being used, the same attack can be applied to the middle DES once the final DES key is known. Finally, if a triple DES cipher with three keys is used, the attack can be iterated to the initial DES.

**AES DFA.** Several DFA attacks have been published for the AES cipher. Faulting a recent smartcard is difficult, with a high risk to cause the chip to self-destruct if it detects an attack; that's why recent DFA research try to minimize the requirements on the number of faults [35,69]. But in the white-box attack model, faults are very easy and cheap to perform. Therefore applying the DFA attack which was introduced in [26] in 2002 is probably one of the best strategies in our white-box context.

We first describe the attack on AES128 encryption, and later discuss how to extend it to AES192 and AES256. For AES128, the final two rounds of the encryption process consist of the following operations: **SubBytes**, **ShiftRows**, **MixColumns**, **AddRoundKey** ( $K_9$ ), **SubBytes**, **ShiftRows**, and **AddRoundKey** ( $K_{10}$ ). Where these operations are defined as follows:

- **SubBytes:** Each byte in the AES state is transformed by applying the AES S-Box.
- **ShiftRows:** The rows in the AES state are shifted to the left by 0, 1, 2 or 3 cells (row 1 to 4).
- **MixColumns:** A matrix multiplication is applied, which results in applying a 32-bit transformation to each column of the state.

We recall how an one-byte fault introduced before the MixColumns in round 9 propagates to the output of the cipher. If the first byte of the state is altered from  $A$  to  $X$ , the fault will propagate to the complete first column after the MixColumns operation.

$$\begin{pmatrix} A & E & I & M \\ B & F & J & N \\ C & G & K & O \\ D & H & L & P \end{pmatrix} \text{ and } \begin{pmatrix} X & E & I & M \\ B & F & J & N \\ C & G & K & O \\ D & H & L & P \end{pmatrix}$$

$$\begin{pmatrix} 2A + 3B + C + D \cdots \cdots \cdots \\ A + 2B + 3C + D \cdots \cdots \cdots \\ A + B + 2C + 3D \cdots \cdots \cdots \\ 3A + B + C + 2D \cdots \cdots \cdots \end{pmatrix} \text{ and } \begin{pmatrix} 2X + 3B + C + D \cdots \cdots \cdots \\ X + 2B + 3C + D \cdots \cdots \cdots \\ X + B + 2C + 3D \cdots \cdots \cdots \\ 3X + B + C + 2D \cdots \cdots \cdots \end{pmatrix}$$

The next steps are performed byte-wise, hence, one faulty byte before the last MixColumns propagates into 4 faulty bytes at the output.

$$\begin{pmatrix} S(2X + 3B + C + D + K_{9,0}) + K_{10,0} & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & S(X + 2B + 3C + D + K_{9,1}) + K_{10,13} \\ \cdots & \cdots & S(X + B + 2C + 3D + K_{9,2}) + K_{10,10} & \cdots \\ \cdots & S(3X + B + C + 2D + K_{9,3}) + K_{10,7} & \cdots & \cdots \end{pmatrix}$$

For the first faulty byte, we can write the following expressions:

$$\begin{aligned} S(2A + 3B + C + D + K_{9,0}) + K_{10,0} &= O_0 \\ S(2X + 3B + C + D + K_{9,0}) + K_{10,0} &= O'_0 \end{aligned}$$

And after xor-ing we obtain

$$S(2A + 3B + C + D + K_{9,0}) \oplus S(2X + 3B + C + D + K_{9,0}) = O_0 \oplus O'_0$$

Similar expressions can be written for each of the faulty bytes, obtaining a set of 4 related equations. Just like in the attack for DES described above, the attack on AES involves solving these equations to obtain a sub-set of candidates for parts of the round key. In this case, each fault provides potential candidates for 4 bytes of the key. Repeating this process with different faults allows finding unique values for each sub-key. Repeating this process for the other columns of the AES state allows recovering candidates for all the sub-keys, and therefore leads to the last round key. For AES128, the last round key is enough to perform a reverse key schedule operation and retrieve the original AES key. As mentioned earlier, the detailed description of this attack can be found in [26]. A variant exists for AES decryption. For AES192 and AES256, two round keys need to be recovered in order to compute the full AES key. Attacking the one before last key requires a small trick to revert the last round in a way to get a succession of operations identical to the original final round: one must permute the *MixColumns* and *AddRoundKey*( $K_{n-1}$ ) steps into *AddRoundKey*(*InvMixColumns*( $K_{n-1}$ )) and *MixColumns*.

## 4.2 Injecting faults into White-box Cryptography

In order to apply DFA attacks to a white-box cryptographic implementation, one fundamental requirement needs to be satisfied: the output needs to be available in a non-encoded form. Once this requirement is satisfied, the attack requires the ability to inject faults into the cryptographic process at the right locations within the algorithm.

In order to locate the appropriate location in which faults need to be injected we typically combine static and dynamic code analysis. For example, recording execution and memory access traces and visualizing them can highlight the location of the target cryptographic algorithm and its rounds, as shown in Fig. 1. We can clearly see the execution of 16 rounds, and therefore we are able to determine at which time and in which memory region faults need to be introduced.

This reverse-engineering step can even be skipped by following another strategy: randomly injecting faults during the computation and observing the output of the cryptographic algorithm can be used to determine the correct location. For example, injecting bit faults in early rounds of a DES execution will result in a fully randomized output. Injecting bit faults during the computation of the last round (i.e. too late for the attack described in Section 4.1) will result in changes to the left half, but not the right half. Similarly, for the AES cipher attack described in Section 4.1, injecting one-byte faults anywhere between the last two *MixColumns* steps will result in 4 corrupted bytes. The location of the 4 corrupted bytes must follow a specific pattern, as indicated by the *MixColumns* and *ShiftRows* combination. Finally, injecting faults can be as simple as flipping bits of the DES or AES intermediate results during the execution of the algorithms. To this end, we can use several techniques:

- If the code can be easily lifted to a high level language (e.g. C/C++), we can introduce code to inject random faults during the target computation. However, in some situations (i.e. with highly obfuscated code) this is a very labour-intensive and error-prone task.
- If the code can be run under a DBI framework (PIN, Valgrind, etc.) we can instrument the code to inject faults and collect the output data.
- A scriptable debugger (e.g. vtrace, gdb) can also be used. To this end, we can write debugger scripts to automate the execution of the cipher, injection of faults and collect the output data.
- Alternatively, emulator-based setups can be used. This allows an attacker to monitor the complete execution and alter the code or data at any desired time. For example, the WBC code can be extracted at runtime and loaded into a standalone emulator such as Unicorn Engine, or a full system emulator such as the PANDA framework could be used.
- In case there is no or only weak self-integrity checks, the tables or even the binary itself could be corrupted to introduce faults statically. A divide and conquer approach can be used by faulting initially large parts of the tables to detect which parts are effectively used when processing a specific input.

Once a way to inject faults is implemented, performing the attack is just a matter of collecting enough faulty output, filter out the non-exploitable ones and plugging the ones with a good fault pattern into the appropriate DFA algorithm.

### 4.3 Analyzing Publicly Available White-Box Implementations

**The Wyseur Challenge.** For the Wyseur challenge, we simply load the ELF binary into a custom emulator script based on the Unicorn Engine. Our custom

emulator records every read and write to memory, and allows to inject faults at selected moments. This approach is fairly trivial for this challenge because the provided binary offers very little functionality and does not interface with many system libraries. The only call to external libraries performed by the binary is used to print the plaintext and ciphertext, which we can simply patch out. We use the execution graph shown in Figure 3a to select the desired interval for the fault injection. We then implement the following process in our emulator:

1. We select a random number within the target interval. We call this the target index.
2. At each memory access, we increment a counter. We then flip a random bit of the first memory read that occurs after the selected target index.
3. At the end of the execution, we record the faulty plaintext.
4. After enough faults are collected, we run the DFA algorithm described in 4.1 to recover the DES key.

We are able to retrieve the key after injecting 40 faults.

**The Hack.lu 2009 Challenge.** For this challenge, we follow a similar approach. The only difference here is that we isolate the white-box code and emulate only the AES function in order to prevent any interaction with the Windows OS. An alternative solution would be to use Intel PIN, or to inject faults in the data section of the binary since no integrity checks are performed. Once the WBC implementation is isolated, the rest of the fault injection process is identical to the Wyseur challenge. The final step is to apply an AES DFA, as described in 4.1. With this attack, we were able to recover the AES key after injecting 90 faults.

**The SSTIC 2012 Challenge.** As described in Section 3.3, the SSTIC challenge was provided as a compiled Python program. We used uncompile2 to lift and recover the original Python code and inject faults into it. The DFA algorithm used was the one presented in 4.1, which allowed us to recover the DES key after injecting 60 faults.

**A White-Box Implementation of the Karroumi Approach.** Since the source code of this implementation is available, we simply injected faults in the original C++ code. We then collected the results and ran the AES DFA attack described in Section 4.1. This allows us to recover the complete AES key after injecting 80 faults.

**The NoSuchCon 2013 Challenge.** As explained in 3.3, the NSC challenge binary incorporates unknown external encodings. For this reason, the effects of the injected faults cannot be easily observed in the output of the white-box. Therefore, the DFA attack described in this paper cannot be directly applied without first recovering the output encoding applied by the implementation. If the output encoding is known or not present, it is possible to retrieve the key in a similar way to the other AES implementations.

#### 4.4 Countermeasures against DFA

Countermeasures against DFA attacks usually involve some sort of redundant computation. For example, a device could perform the encryption process twice and compare both outputs. Assuming that the adversary does not have any access to the intermediate data, it is possible to detect the attack and prevent outputting faulty results. However, in the white-box settings this direct approach is not valid: if the attacker is able to observe the comparison of the two results, he can simply duplicate the faulty result and bypass the check. In order to protect a white-box cryptographic implementation against DFA, the following two avenues might be used:

- Carrying redundant data (e.g. error-correcting codes) along with the computation in such a way that a modification performed by an attacker can be compensated, without affecting the data in the non-encoded domain.
- Implementing the internal data encodings in such a way that faults propagate into larger parts of the cipher state. In this way, the fault models expected by the standard DFA attacks do not apply and therefore an attacker would have to develop customized attacks.

## 5 Conclusions and Future Work

As conjectured in the first papers introducing the white-box attack model, one cannot expect long-term defense against attacks on white-box implementations. However, as we have shown in this work, all current publicly available white-box implementations do not even offer any short-term security since the DCA and DFA techniques can extract the secret key within seconds.

Although we sketched some ideas on countermeasures, it remains an open question how to guard oneself against these types of attacks. The countermeasures against differential power analysis attacks applied in the area of high-assurance applications do not seem to carry over directly due to the ability of the adversary to disable or tamper with the random source. If medium to long term security is required then tamper resistant hardware solutions, like a secure element, seem like a much better alternative.

Similarly, all publicly known differential fault analysis countermeasures are based on introducing redundancy in the computation. If this is not done carefully, an attacker might still be able to inject faults while bypassing the redundancy checks. Additional research on improving the robustness of WBC implementations against fault attacks is thus required.

Another interesting research direction is to see if the more advanced and powerful techniques used in side-channel analysis from the cryptographic hardware community obtain even better results in this setting. Examples include correlation power analysis and higher order attacks. Studying other point of attack, for instance targeting the multiplicative inverse step in the first round of AES, give interesting results (see Section 3.3). Investigating other positions as a target in our DCA approach may be worth as well.

## References

1. Advanced Encryption Standard (AES). National Institute of Standards and Technology (NIST), FIPS PUB 197, U.S. Department of Commerce, Nov. 2001.
2. B. Amstadt and M. K. Johnson. Wine. *Linux Journal*, 1994(4), August 1994.
3. C. H. Baek, J. H. Cheon, and H. Hong. Analytic toolbox for white-box implementations: Limitation and perspectives. Cryptology ePrint Archive, Report 2014/688, 2014. <http://eprint.iacr.org/2014/688>.
4. B. Barak, S. Garg, Y. T. Kalai, O. Paneth, and A. Sahai. Protecting obfuscation against algebraic attacks. In P. Q. Nguyen and E. Oswald, editors, *EUROCRYPT 2014*, volume 8441 of *LNCS*, pages 221–238, Copenhagen, Denmark, May 11–15, 2014. Springer, Heidelberg, Germany.
5. B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. P. Vadhan, and K. Yang. On the (im)possibility of obfuscating programs. In J. Kilian, editor, *CRYPTO 2001*, volume 2139 of *LNCS*, pages 1–18, Santa Barbara, CA, USA, Aug. 19–23, 2001. Springer, Heidelberg, Germany.
6. J.-B. Bédrupe. Hack.lu 2009 reverse challenge 1. online, 2009. <http://2009.hack.lu/index.php/ReverseChallenge>.
7. F. Bellard. QEMU, a fast and portable dynamic translator. In *USENIX Annual Technical Conference, FREENIX Track*, pages 41–46, 2005.
8. S. Bhatkar, D. C. DuVarney, and R. Sekar. Address obfuscation: An efficient approach to combat a broad range of memory error exploits. In *Proceedings of the 12th USENIX Security Symposium*. USENIX Association, 2003.
9. E. Biham and A. Shamir. Differential cryptanalysis of Snefru, Khafre, REDOC-II, LOKI and Lucifer. In J. Feigenbaum, editor, *CRYPTO'91*, volume 576 of *LNCS*, pages 156–171, Santa Barbara, CA, USA, Aug. 11–15, 1992. Springer, Heidelberg, Germany.
10. E. Biham and A. Shamir. Differential fault analysis of secret key cryptosystems. In B. S. Kaliski Jr., editor, *CRYPTO'97*, volume 1294 of *LNCS*, pages 513–525, Santa Barbara, CA, USA, Aug. 17–21, 1997. Springer, Heidelberg, Germany.
11. O. Billet and H. Gilbert. A traceable block cipher. In C.-S. Lai, editor, *ASIACRYPT 2003*, volume 2894 of *LNCS*, pages 331–346. Springer, Heidelberg, Germany, 2003.
12. O. Billet, H. Gilbert, and C. Ech-Chatbi. Cryptanalysis of a white box AES implementation. In H. Handschuh and A. Hasan, editors, *SAC 2004*, volume 3357 of *LNCS*, pages 227–240, Waterloo, Ontario, Canada, Aug. 9–10, 2004. Springer, Heidelberg, Germany.
13. A. Biryukov, C. Bouillaguet, and D. Khovratovich. Cryptographic schemes based on the ASASA structure: Black-box, white-box, and public-key (extended abstract). In P. Sarkar and T. Iwata, editors, *ASIACRYPT 2014, Part I*, volume 8873 of *LNCS*, pages 63–84, Kaoshiung, Taiwan, R.O.C., Dec. 7–11, 2014. Springer, Heidelberg, Germany.
14. A. Biryukov, C. De Cannière, A. Braeken, and B. Preneel. A toolbox for cryptanalysis: Linear and affine equivalence algorithms. In E. Biham, editor, *EUROCRYPT 2003*, volume 2656 of *LNCS*, pages 33–50, Warsaw, Poland, May 4–8, 2003. Springer, Heidelberg, Germany.
15. D. Boneh, R. A. DeMillo, and R. J. Lipton. On the importance of checking cryptographic protocols for faults (extended abstract). In W. Fumy, editor, *EUROCRYPT'97*, volume 1233 of *LNCS*, pages 37–51, Konstanz, Germany, May 11–15, 1997. Springer, Heidelberg, Germany.

16. Z. Brakerski and G. N. Rothblum. Virtual black-box obfuscation for all circuits via generic graded encoding. In Y. Lindell, editor, *TCC 2014*, volume 8349 of *LNCS*, pages 1–25, San Diego, CA, USA, Feb. 24–26, 2014. Springer, Heidelberg, Germany.
17. E. Brier, C. Clavier, and F. Olivier. Correlation power analysis with a leakage model. In M. Joye and J.-J. Quisquater, editors, *CHES 2004*, volume 3156 of *LNCS*, pages 16–29, Cambridge, Massachusetts, USA, Aug. 11–13, 2004. Springer, Heidelberg, Germany.
18. J. Bringer, H. Chabanne, and E. Dottax. White box cryptography: Another attempt. Cryptology ePrint Archive, Report 2006/468, 2006. <http://eprint.iacr.org/2006/468>.
19. S. Chari, C. S. Jutla, J. R. Rao, and P. Rohatgi. Towards sound approaches to counteract power-analysis attacks. In M. J. Wiener, editor, *CRYPTO'99*, volume 1666 of *LNCS*, pages 398–412, Santa Barbara, CA, USA, Aug. 15–19, 1999. Springer, Heidelberg, Germany.
20. S. Chari, J. R. Rao, and P. Rohatgi. Template attacks. In B. S. Kaliski Jr., Çetin Kaya. Koç, and C. Paar, editors, *CHES 2002*, volume 2523 of *LNCS*, pages 13–28, Redwood Shores, CA, USA, Aug. 13–15, 2003. Springer, Heidelberg, Germany.
21. S. Chow, P. A. Eisen, H. Johnson, and P. C. van Oorschot. White-box cryptography and an AES implementation. In K. Nyberg and H. M. Heys, editors, *SAC 2002*, volume 2595 of *LNCS*, pages 250–270, St. John's, Newfoundland, Canada, Aug. 15–16, 2003. Springer, Heidelberg, Germany.
22. S. Chow, P. A. Eisen, H. Johnson, and P. C. van Oorschot. A white-box DES implementation for DRM applications. In J. Feigenbaum, editor, *Security and Privacy in Digital Rights Management, ACM CCS-9 Workshop, DRM 2002*, volume 2696 of *LNCS*, pages 1–15. Springer, 2003.
23. J. Daemen and V. Rijmen. *The design of Rijndael: AES — the Advanced Encryption Standard*. Springer, 2002.
24. Y. de Mulder. *White-Box Cryptography: Analysis of White-Box AES Implementations*. PhD thesis, KU Leuven, 2014.
25. C. Delerablée, T. Lepoint, P. Paillier, and M. Rivain. White-box security notions for symmetric encryption schemes. In T. Lange, K. Lauter, and P. Lisonek, editors, *SAC 2013*, volume 8282 of *LNCS*, pages 247–264, Burnaby, BC, Canada, Aug. 14–16, 2014. Springer, Heidelberg, Germany.
26. P. Dusart, G. Letourneux, and O. Vivolo. Differential fault analysis on A.E.S. In J. Zhou, M. Yung, and Y. Han, editors, *ACNS 2003*, volume 2846 of *Lecture Notes in Computer Science*, pages 293–306. Springer, 2003.
27. F. Falco and N. Riva. Dynamic binary instrumentation frameworks: I know you're there spying on me. REcon, 2012. <http://recon.cx/2012/schedule/events/216.en.html>.
28. S. Garg, C. Gentry, S. Halevi, M. Raykova, A. Sahai, and B. Waters. Candidate indistinguishability obfuscation and functional encryption for all circuits. In *54th Annual IEEE Symposium on Foundations of Computer Science, FOCS*, pages 40–49. IEEE Computer Society, 2013.
29. L. Goubin, J.-M. Masereel, and M. Quisquater. Cryptanalysis of white box DES implementations. In C. M. Adams, A. Miri, and M. J. Wiener, editors, *SAC 2007*, volume 4876 of *LNCS*, pages 278–295, Ottawa, Canada, Aug. 16–17, 2007. Springer, Heidelberg, Germany.
30. L. Goubin and J. Patarin. DES and differential power analysis (the “duplication” method). In Çetin Kaya. Koç and C. Paar, editors, *CHES'99*, volume 1717 of

- LNCS*, pages 158–172, Worcester, Massachusetts, USA, Aug. 12–13, 1999. Springer, Heidelberg, Germany.
31. Y. Huang, F. S. Ho, H. Tsai, and H. M. Kao. A control flow obfuscation method to discourage malicious tampering of software codes. In F. Lin, D. Lee, B. P. Lin, S. Shieh, and S. Jajodia, editors, *Proceedings of the 2006 ACM Symposium on Information, Computer and Communications Security, ASIACCS 2006*, page 362. ACM, 2006.
  32. M. Jacob, D. Boneh, and E. W. Felten. Attacking an obfuscated cipher by injecting faults. In J. Feigenbaum, editor, *Security and Privacy in Digital Rights Management, ACM CCS-9 Workshop, DRM 2002, Washington, DC, USA, November 18, 2002, Revised Papers*, volume 2696 of *LNCS*, pages 16–31. Springer, 2003.
  33. M. Jakobsson and M. K. Reiter. Discouraging software piracy using software aging. In T. Sander, editor, *Security and Privacy in Digital Rights Management, ACM CCS-8 Workshop DRM 2001*, volume 2320 of *LNCS*, pages 1–12. Springer, 2002.
  34. M. Karroumi. Protecting white-box AES with dual ciphers. In K. H. Rhee and D. Nyang, editors, *ICISC 10*, volume 6829 of *LNCS*, pages 278–291, Seoul, Korea, Dec. 1–3, 2011. Springer, Heidelberg, Germany.
  35. C. H. Kim and J. Quisquater. New differential fault analysis on AES key schedule: Two faults are enough. In G. Grimaud and F. Standaert, editors, *CARDIS 2008*, volume 5189 of *Lecture Notes in Computer Science*, pages 48–60. Springer, 2008.
  36. J. Kirsch. Towards transparent dynamic binary instrumentation using virtual machine introspection. REcon, 2015. <https://recon.cx/2015/schedule/events/20.html>.
  37. D. Klinec. White-box attack resistant cryptography. Master’s thesis, Masaryk University, Brno, Czech Republic, 2013. [https://is.muni.cz/th/325219/fi\\_m/](https://is.muni.cz/th/325219/fi_m/).
  38. P. Kocher, J. Jaffe, B. Jun, and P. Rohatgi. Introduction to differential power analysis. *Journal of Cryptographic Engineering*, 1(1):5–27, 2011.
  39. P. C. Kocher, J. Jaffe, and B. Jun. Differential power analysis. In M. J. Wiener, editor, *CRYPTO’99*, volume 1666 of *LNCS*, pages 388–397, Santa Barbara, CA, USA, Aug. 15–19, 1999. Springer, Heidelberg, Germany.
  40. T. Lepoint, M. Rivain, Y. D. Mulder, P. Roelse, and B. Preneel. Two attacks on a white-box AES implementation. In T. Lange, K. Lauter, and P. Lisonek, editors, *SAC 2013*, volume 8282 of *LNCS*, pages 265–285, Burnaby, BC, Canada, Aug. 14–16, 2014. Springer, Heidelberg, Germany.
  41. X. Li and K. Li. Defeating the transparency features of dynamic binary instrumentation. BlackHat US, 2014. <https://www.blackhat.com/docs/us-14/materials/us-14-Li-Defeating-The-Transparency-Feature-Of-DBI.pdf>.
  42. H. E. Link and W. D. Neumann. Clarifying obfuscation: Improving the security of white-box DES. In *International Symposium on Information Technology: Coding and Computing (ITCC 2005)*, pages 679–684. IEEE Computer Society, 2005.
  43. C. Linn and S. K. Debray. Obfuscation of executable code to improve resistance to static disassembly. In S. Jajodia, V. Atluri, and T. Jaeger, editors, *Proceedings of the 10th ACM Conference on Computer and Communications Security, CCS 2003*, pages 290–299. ACM, 2003.
  44. C. Luk, R. S. Cohn, R. Muth, H. Patil, A. Klauser, P. G. Lowney, S. Wallace, V. J. Reddi, and K. M. Hazelwood. Pin: building customized program analysis tools with dynamic instrumentation. In V. Sarkar and M. W. Hall, editors, *Proceedings of the ACM SIGPLAN 2005 Conference on Programming Language Design and Implementation*, pages 190–200. ACM, 2005.

45. A. Maillet. Nosuchcon 2013 challenge - write up and methodology. online, 2013. <http://kutioo.blogspot.be/2013/05/nosuchcon-2013-challenge-write-up-and.html>.
46. S. Mangard, E. Oswald, and F. Standaert. One for all - all for one: unifying standard differential power analysis attacks. *IET Information Security*, 5(2):100–110, 2011.
47. F. Marceau, F. Perigaud, and A. Tillequin. Challenge SSTIC 2012. online, 2012. <http://communaute.sstic.org/ChallengeSSTIC2012>.
48. T. S. Messerges. Using second-order power analysis to attack DPA resistant software. In Çetin Kaya. Koç and C. Paar, editors, *CHES 2000*, volume 1965 of *LNCS*, pages 238–251, Worcester, Massachusetts, USA, Aug. 17–18, 2000. Springer, Heidelberg, Germany.
49. W. Michiels. Opportunities in white-box cryptography. *IEEE Security & Privacy*, 8(1):64–67, 2010.
50. W. Michiels and P. Gorissen. Mechanism for software tamper resistance: an application of white-box cryptography. In M. Yung, A. Kiayias, and A. Sadeghi, editors, *Proceedings of the Seventh ACM Workshop on Digital Rights Management*, pages 82–89. ACM, 2007.
51. W. Michiels, P. Gorissen, and H. D. L. Hollmann. Cryptanalysis of a generic class of white-box implementations. In R. M. Avanzi, L. Keliher, and F. Sica, editors, *SAC 2008*, volume 5381 of *LNCS*, pages 414–428, Sackville, New Brunswick, Canada, Aug. 14–15, 2009. Springer, Heidelberg, Germany.
52. C. Mougey and F. Gabriel. Désobfuscation de DRM par attaques auxiliaires. In *Symposium sur la sécurité des technologies de l'information et des communications*, 2014. [www.sstic.org/2014/presentation/dsobfuscation\\_de\\_drm\\_par\\_attaques\\_auxiliaires](http://www.sstic.org/2014/presentation/dsobfuscation_de_drm_par_attaques_auxiliaires).
53. J. A. Muir. A tutorial on white-box AES. In E. Kranakis, editor, *Advances in Network Analysis and its Applications*, volume 18 of *Mathematics in Industry*, pages 209–229. Springer Berlin Heidelberg, 2013.
54. Y. D. Mulder, P. Roelse, and B. Preneel. Cryptanalysis of the Xiao-Lai white-box AES implementation. In L. R. Knudsen and H. Wu, editors, *SAC 2012*, volume 7707 of *LNCS*, pages 34–49, Windsor, Ontario, Canada, Aug. 15–16, 2013. Springer, Heidelberg, Germany.
55. Y. D. Mulder, B. Wyseur, and B. Preneel. Cryptanalysis of a perturbed white-box AES implementation. In G. Gong and K. C. Gupta, editors, *INDOCRYPT 2010*, volume 6498 of *LNCS*, pages 292–310, Hyderabad, India, Dec. 12–15, 2010. Springer, Heidelberg, Germany.
56. N. Nethercote and J. Seward. Valgrind: a framework for heavyweight dynamic binary instrumentation. In J. Ferrante and K. S. McKinley, editors, *Proceedings of the ACM SIGPLAN 2007 Conference on Programming Language Design and Implementation*, pages 89–100. ACM, 2007.
57. S. Nikova, C. Rechberger, and V. Rijmen. Threshold implementations against side-channel attacks and glitches. In P. Ning, S. Qing, and N. Li, editors, *Information and Communications Security, ICICS*, volume 4307 of *LNCS*, pages 529–545. Springer, 2006.
58. J. Patarin and L. Goubin. Asymmetric cryptography with S-boxes. In Y. Han, T. Okamoto, and S. Qing, editors, *ICICS 97*, volume 1334 of *LNCS*, pages 369–380, Beijing, China, Nov. 11–14, 1997. Springer, Heidelberg, Germany.
59. M. L. Polla, F. Martinelli, and D. Sgandurra. A survey on security for mobile devices. *IEEE Communications Surveys and Tutorials*, 15(1):446–471, 2013.

60. P. Sasdrich, A. Moradi, and T. Güneysu. White-box cryptography in the gray box - - A hardware implementation and its side channels -. In T. Peyrin, editor, *FSE 2016*, volume 9783 of *LNCS*, pages 185–203, Bochum, Germany, Mar. 20–23, 2016. Springer, Heidelberg, Germany.
61. A. Saxena, B. Wyseur, and B. Preneel. Towards security notions for white-box cryptography. In P. Samarati, M. Yung, F. Martinelli, and C. A. Ardagna, editors, *ISC 2009*, volume 5735 of *LNCS*, pages 49–58, Pisa, Italy, Sept. 7–9, 2009. Springer, Heidelberg, Germany.
62. K. Schramm and C. Paar. Higher order masking of the AES. In D. Pointcheval, editor, *CT-RSA 2006*, volume 3860 of *LNCS*, pages 208–225, San Jose, CA, USA, Feb. 13–17, 2006. Springer, Heidelberg, Germany.
63. F. Scrinzi. Behavioral analysis of obfuscated code. Master’s thesis, University of Twente, Twente, Netherlands, 2015. [http://essay.utwente.nl/67522/1/Scrinzi\\_MA\\_SCS.pdf](http://essay.utwente.nl/67522/1/Scrinzi_MA_SCS.pdf).
64. A. Souchet. AES whitebox unboxing: No such problem. online, 2013. <http://overclock.tuxfamily.org/blog/?p=253>.
65. SysK. Practical cracking of white-box implementations. Phrack 68:14. <http://www.phrack.org/issues/68/8.html>.
66. P. Teuwen. CHES2015 writeup. online, 2015. [http://wiki.yobi.be/wiki/CHES2015\\_Writeup#Challenge\\_4](http://wiki.yobi.be/wiki/CHES2015_Writeup#Challenge_4).
67. P. Teuwen. NSC writeups. online, 2015. [http://wiki.yobi.be/wiki/NSC\\_Writeups](http://wiki.yobi.be/wiki/NSC_Writeups).
68. L. Tolhuizen. Improved cryptanalysis of an AES implementation. In *Proceedings of the 33rd WIC Symposium on Information Theory*. Werkgemeenschap voor Inform-en Communicatietheorie, 2012.
69. M. Tunstall, D. Mukhopadhyay, and S. Ali. Differential fault analysis of the advanced encryption standard using a single fault. In C. A. Ardagna and J. Zhou, editors, *WISTP 2011*, volume 6633 of *Lecture Notes in Computer Science*, pages 224–233. Springer, 2011.
70. U.S. DEPARTMENT OF COMMERCE/National Institute of Standards and Technology. *Data Encryption Standard (DES)*.
71. E. Vanderbéken. Hacklu reverse challenge write-up. online, 2009. <http://baboon.rce.free.fr/index.php?post/2009/11/20/HackLu-Reverse-Challenge>.
72. B. Wyseur, W. Michiels, P. Gorissen, and B. Preneel. Cryptanalysis of white-box DES implementations with arbitrary external encodings. In C. M. Adams, A. Miri, and M. J. Wiener, editors, *SAC 2007*, volume 4876 of *LNCS*, pages 264–277, Ottawa, Canada, Aug. 16–17, 2007. Springer, Heidelberg, Germany.
73. Y. Xiao and X. Lai. A secure implementation of white-box AES. In *Computer Science and its Applications, 2009. CSA '09. 2nd International Conference on*, pages 1–6, 2009.
74. Y. Zhou and S. Chow. System and method of hiding cryptographic private keys, Dec. 15 2009. US Patent 7,634,091.