

Attacking an obfuscated cipher by injecting faults

Matthias Jacob
mjacob@cs.princeton.edu

Dan Boneh
dabo@cs.stanford.edu

Edward Felten
felten@cs.princeton.edu

Abstract

We study the strength of certain obfuscation techniques used to protect software from reverse engineering and tampering. We show that some common obfuscation methods can be defeated using a fault injection attack, namely an attack where during program execution an attacker injects errors into the program environment. By observing how the program fails under certain errors the attacker can deduce the obfuscated information in the program code without having to unravel the obfuscation mechanism. We apply this technique to extract a secret key from a block cipher obfuscated using a commercial obfuscation tool and draw conclusions on preventing this weakness.

1 Introduction

In recent years the advent of mass distribution of digital content fueled the demand for tools to prevent software and digital media from illegal copying. The goal is to make it harder for a malicious person to reverse engineer or modify a given piece of software. One well known technique for preventing illegal use of digital media is watermarking for audio and video content [23] which had only limited success. Another common approach is to only distribute encrypted content (see, e.g., CSS [2], Intertrust [3], MS Windows Media Technologies [4], Adobe EBooks [1]). Users run content players on their machines and these players enforce access permissions associated with the content. In most of these systems the software player contains some secret information that enables it to decrypt the content internally. Clearly the whole point is that the user should not be able to emulate the player and decrypt the content by herself. As a result, the secret information that enables the player to decrypt the content must be hidden somehow in the player's binary code. We note that hardware solutions, where the decryption key is embedded in tamper-resistant hardware [13, 7, 6], have had some success [14, 28], but clearly a software only solution, assuming it is secure, is superior because it is more cost efficient and easier to deploy.

This brings us to one of the main challenges facing content protection vendors: is it possible to hide a decryption key in the implementation of a block cipher (e.g. AES) in such a way that given the binary code it is hard to extract the decryption key. In other words, suppose $D^k(c)$ is an algorithm for decrypting the ciphertext c using the key k . Is it possible to modify the implementation of $D^k(c)$ so that extracting k by reverse engineering is sufficiently hard? If hiding the key in a binary is possible, it has a crucial advantage over alternative key hiding techniques: in order to decrypt content the binary needs to be executed, and efficient access control mechanisms exist in the operating system in order to prevent unauthorized execution, whereas hiding a stored key in memory is difficult [33]. Key

obfuscation is a very old question already mentioned in the classic paper of Diffie and Hellman [25].

Code obfuscation is a common technique for protecting software against reverse engineering and is commonly used for hiding proprietary software systems and sensitive system components such as a cipher. Commercial obfuscation tools often work by taking as input arbitrary program source code, and they output obfuscated binary or source code that is harder to reverse engineer and thus to manipulate than the original software [12, 8, 10, 9, 5]. However, it is unclear whether obfuscation techniques can be strong enough to protect sensitive software systems such as a cipher implementation.

In this paper we investigate a commercial state-of-the-art obfuscated cryptosystem [21] that hides a secret key. An *ideal obfuscation tool* turns program code into a black-box, and therefore it is impossible to find out any properties of the program. In practice however, obfuscation tools often only *approximate the ideal case*. When obfuscating a cryptosystem the obfuscator embeds a secret key into the program code and obfuscates the code. It should be hard to figure out any properties about the key by just investigating the code. However, we show how to extract the secret key from the system in only a few cryptographic operations and come to the conclusion that current obfuscation techniques for hiding a secret key are not strong enough to resist certain attacks.

Our attack is based on differential fault analysis [17] in which an attacker injects errors into the code in order to get information about the secret key. The impact of this attack is comparable to an attack on an RSA implementation based on the Chinese Remainder Theorem that requires only one faulty RSA signature in order to extract the private key [18].

Fault attacks are a threat on tamper-resistant hardware [14], and in this paper we show that an adversary can also inject faults to extract a key from obfuscated software. Based on our experience in attacking an obfuscated cryptosystem we propose techniques for strengthening code obfuscation to make fault attacks more difficult and make a first step in understanding the limits of practical software obfuscation.

2 Attacking an obfuscated cipher implementation

In this section we describe our attack on a state-of-the-art obfuscator [21] illustrated in Figure 1. We were given the obfuscated source code for both DES encryption and decryption of the iterated block cipher. Our goal was to reverse engineer the system only based on knowledge of this obfuscated source code. For the given obfuscated code the attacker does not learn more properties about the program by investigating the obfuscated source code than by just disassembling the binary because most of the program is composed of lookup tables.

In this particular approach the obfuscation method hides the secret key of a round-based cipher in the code. Because a round-based cipher exposes the secret key every time it combines the key with the input data of a round, the obfuscator injects randomness and redundancies and refines the resulting boolean operations into lookup tables. Instead of executing algorithmic code, the program steps through a chain of precomputed values in lookup tables and retrieves the correct result. Therefore it is difficult to obtain any information about the single rounds by just looking at the source code or binary code, but in our attack we obtain information by observing and changing data during the encryption process.

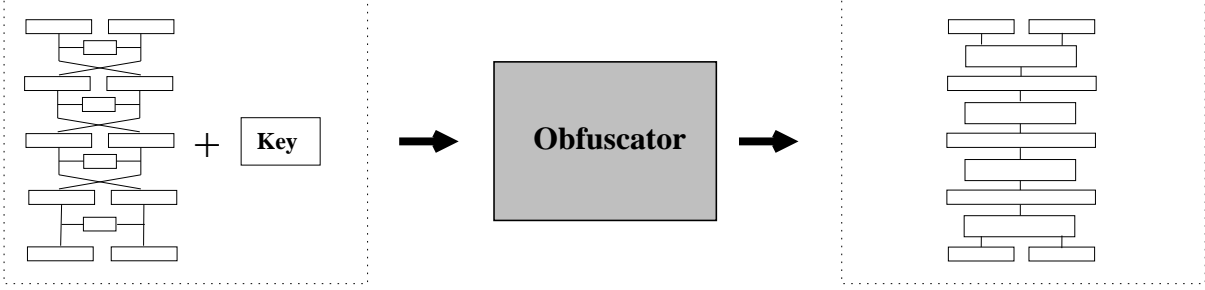


Figure 1: Operation of the obfuscator on the round-based cipher: It transforms the key and the original source code into code that implements every round as a lookup table of precomputed values. The intermediate results after each round are encoded.

2.1 Obfuscating an iterated block cipher

The obfuscation process of the cipher implementation is shown in Figure 1. The obfuscator transforms the original source code and the key into a cipher in which the key is embedded and hidden in the rounds. The single rounds of the cipher are unrolled, but the boundaries of each round are clearly recognizable. The cipher contains n rounds π_i^k for each $i = 1, \dots, n$ with the key k . Including the initial permutation λ the cipher computes the function

$$E^k(M) := \left[\lambda^{-1} \cdot \pi_n^k \cdot \pi_{n-1}^k \cdot \dots \cdot \pi_1^k \cdot \lambda \right] (M).$$

However, interpretation of any intercepted intermediate results is difficult since the obfuscator maps the original intermediate results after each round to a new representation. This transformation is described in detail in [21].

In the following paragraphs we give an algebraic definition for the transformation into the 96-bit intermediate representation of the obfuscator in [21]. In the first step we define some basic operations. $x|_i^m$ extracts bits i through $i + m$ from a bit string. $EP(x)$ computes the DES expansion permutation.

$$\begin{aligned} x_1x_2\dots x_n|_i^m &= x_ix_{i+1}\dots x_{i+m} \\ x_1x_2\dots x_n|_i &= x_i \end{aligned}$$

$$EP_i(x) = EP(x)|_{6i}^6$$

$$\begin{aligned} R_r^k &= EP(R_r^k) \\ R_{r,i}^k &= EP_i(R_r^k) \end{aligned}$$

The t-box $T_{r,i}^k(L_r, R_r^k)$ computes the i -th DES s-box in round r for $i = 0..7$ and appends $R(L_r, R_r^k)$ which takes the first and sixth bit from $R_{r,i}^k$ and appends two random bits from L_r . The bits from L_r are used to forward the left hand side information in the t-boxes, and the first and sixth bit from $R_{r,i}^k$ to reconstruct R_r^k from the s-box result in order to forward it to round $r + 1$ as the left hand side input.

$$T_{r,i}^k(L_r, R_r^k) = S_{r,i}^k(R_{r,i}^k) \parallel R(L_r, R_{r,i}^k)$$

$$T_r^k(L_r, R_r^k) = T_{r,\gamma_r(0)}^k(L_r, R_r^k) \parallel T_{r,\gamma_r(1)}^k(L_r, R_r^k) \parallel \dots \parallel T_{r,\gamma_r(11)}^k(L_r, R_r^k)$$

For $i = 8..11$ $T_{r,i}^k(L_r, R_r^k)$ outputs either random dummy values or bits from L_r .

In order to obfuscate the result γ_r permutes the order of the t-boxes on $T_r = \{T_{r,0}^k \dots T_{r,11}^k\}$. Additionally, ϕ_r applies a bijective non-linear encoding on 4-bit blocks x_j for $j = 1..24$ where

$\phi_r(x) = (\phi_{r,1}(x_1), \phi_{r,2}(x_2), \dots, \phi_{r,24}(x_{24}))$ and $x = x_1x_2\dots x_{24}$. Since a single t-box consists of 8 bit outputs, two different bijective non-linear encodings belong to one t-box.

In order to do the second step the obfuscated DES implementation needs to be able to recover the original right hand side input to round r , and this gets implemented using function $\alpha_{r,i}^k(y)$ which takes the forwarded bits x_1 and x_2 that describe the row of the s-box.

$$\alpha_{r,i}^k(y, x_1, x_2) = EP_i^{-1}((S_{r,i}^k)^{-1}(y, x_1, x_2))$$

$$\begin{aligned} L_r &= L_r^0 \parallel L_r^1 \parallel L_r^2 \parallel \dots \parallel L_r^7 \\ R'_r &= R_r'^0 \parallel R_r'^1 \parallel R_r'^2 \parallel \dots \parallel R_r'^7 \end{aligned}$$

The second step then implements the function $\tau_{r,i}^k$ in which $\mu_r(n)$ describes the corresponding position of the bit in the output of the t-boxes, and PB is the DES p-box operation:

$$\begin{aligned} \tau_{r,i}^k(x)(L_r^i, R_r'^i) &= \left(\underbrace{\alpha_{r,i}^k(x|_{8\gamma_r(i)}, x|_{8\gamma_r(i)+4}, x|_{8\gamma_r(i)+5})}_{\text{depends on } R_{r-1} \text{ only}} \parallel \right. \\ &\quad \left. \left(EP_i \left[PB \left(\underbrace{x|_{\gamma_r(0)}^4 \parallel x|_{\gamma_r(1)}^4 \parallel \dots \parallel x|_{\gamma_r(11)}^4}_{\text{depends on } R_{r-1} \text{ only}} \right) \oplus \left(\underbrace{x|_{\mu_r(0)} \parallel \dots \parallel x|_{\mu_r(32)}}_{\text{depends on } L_{r-1} \text{ only}} \right) \right] \right) \right) \\ \tau_r^k(x) &= \tau_{r,0}^k(x) \parallel \tau_{r,1}^k(x) \parallel \dots \parallel \tau_{r,11}^k(x) \end{aligned}$$

ψ_r and ϕ_r are different non-linear bijective encodings on 4-bit blocks, and δ_r

$$\delta_r(L, R') = \gamma_r(\mu_r((L|0^{24}), R'))$$

$$\begin{aligned} \mu_r(x_0x_1\dots x_{47}, y_0\dots y_{47}) &= y_0\dots y_5x_{\mu_r^{-1}(0)}x_{\mu_r^{-1}(1)}y_6\dots y_{11}x_{\mu_r^{-1}(2)}x_{\mu_r^{-1}(3)}\dots y_{42}\dots y_{47}x_{\mu_r^{-1}(22)}x_{\mu_r^{-1}(23)}\dots x_{\mu_r^{-1}(47)} \\ \gamma_r(z_0z_1\dots z_{95}) &= z_{\gamma_r^{-1}(0)}\dots z_{(\gamma_r^{-1}(0)+5)}z_6z_7\dots z_{\gamma_r^{-1}(11)}\dots z_{(\gamma_r^{-1}(11)+5)}z_{94}z_{95} \end{aligned}$$

The obfuscated t-box is

$$T_r'^k(x) = (\phi_r T_r^k \psi_{r-1}^{-1})(x).$$

Hence the transformed function is:

$$E^k(x) = \left[(\lambda^{-1} \delta_n^{-1} \psi_n^{-1}) \cdot (\psi_n \delta_n \tau_n^k \phi_n^{-1}) \cdot (\phi_n T_n^k \psi_{n-1}^{-1}) \cdot \dots \cdot (\psi_1 \delta_1 \tau_1^k \phi_1^{-1}) \cdot (\phi_1 T_1^k \psi_0^{-1}) \cdot (\psi_0 \delta_0 \beta \lambda) \right] (x)$$

with

$$\beta(L, R) = L \parallel EP(R)$$

By setting

$$\tau_r'^k = \begin{cases} \psi_0 \delta_0 \beta \lambda & r = 0 \\ \psi_r \delta_r \tau_r^k \phi_r^{-1} & r = 1, \dots, n \\ \lambda^{-1} \delta_n^{-1} \psi_n^{-1} & r = n + 1 \end{cases}$$

the resulting encryption operation is

$$E^k(x) = \left[\tau_{n+1}'^k \cdot (\tau_n'^k \cdot T_n^k) \cdot \dots \cdot (\tau_1'^k \cdot T_1^k) \cdot \tau_0'^k \right] (x)$$

Every component $\tau_i'^k$ and T_i^k is implemented within a separate lookup table.

For convenience set

$$\tau_r''^k = \begin{cases} \tau_r'^k & r = 0, r = n + 1 \\ \tau_r'^k \cdot T_r'^k & r = 1, \dots, n \end{cases}$$

and obtain

$$E^k(x) = \left[\tau_{n+1}''^k \cdot \tau_n''^k \cdot \dots \cdot \tau_0''^k \right] (x)$$

Figure 2 shows the deobfuscation problem. Given one DES round and the obfuscated intermediate representations an attacker wants to find out the intermediate representation which is encoded by the unknown function σ_r . This σ_r is the inverse of the encoded input to the t-box (by ψ), the permutation of the t-boxes γ_r , and the random distribution of the left hand side μ_r :

$$\sigma_r(L_r, R_r) = \psi_r(\delta_r(L_r, EP(R_r)))$$

$E^k(x)$ contains the key k implicitly in $\tau_r''^k$ (in [21] $\tau_0'^k$ corresponds to M_1 , τ_{n+1}^k to M_3 and all other $\tau_r'^k$ to M_2). In other words, the implementation of $\tau_r''^k$ hides the decomposition into its components σ_{r-1}^{-1} , π_r^k , and σ_r . Hence, recovering the key boils down to the problem of extracting π_r^k out of $\tau_r''^k$. In any further explanations we remove λ from any computation since it does not play any role in the attack and can be easily inverted. Therefore $\tau_0''^k = \psi_0$ and $\tau_{n+1}''^k = \psi_n$.

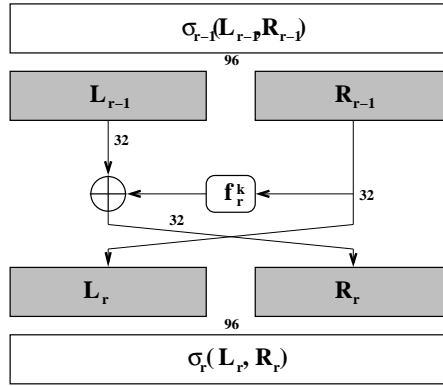


Figure 2: Round r with the function f_r^k hiding the key k . σ_r is the intermediate representation and L_r and R_r are the left hand and the right hand side of the intermediate result respectively. The rounds π_r^k correspond to $\pi_r^k = f_r^k(R_{r-1} \oplus L_{r-1}, R_{r-1})$ for $r = 1..n$.

2.2 Attacking an obfuscated iterated block cipher

In an example for a naive approach for attacking the obfuscated cipher an adversary encrypts some arbitrary plaintext and intercepts intermediate results to obtain $\sigma_r(L_r, R_r)$. The adversary starts the attack by encrypting plaintexts p that have one single bit set, and afterward examines the obfuscated intermediate results after the first round π_1^k during encryption. By heuristically computing the differences between $(\tau_1''\tau_0'')(p)$ and $(\tau_1''\tau_0'')(0)$ for $p \neq 0$ we find that $(\tau_1''\tau_0'')(p)$ changes deterministically for all p that have one bit set in the left hand side of the plaintext L_0 due to the construction of the t-boxes. However, since the adversary is not able to compute σ_1^{-1} in order to retrieve R_1 any knowledge of R_0 and L_0 is meaningless if she wants to extract the key. An attack that works on the first round by recovering σ_1^{-1} of the cipher is the statistical bucketing attack [21]. This attack exploits some properties of the DES s-boxes and requires about 2^{13} encryptions. In contrast our attack works for any round-based block cipher and requires only dozens of encryptions.

We now describe how we use a simplified differential cryptanalysis called differential fault analysis [17] to recover the key in a few operations. In this attack an adversary flips bits in the input to the last round function f_n^k and computes the different outputs to find out the round function f_n^k of the last round n . When injecting single bit faults into the last round using chosen ciphertexts only dozens of cryptographic operations are necessary in order to find f_n^k . The implementation of this attack requires less information about the intermediate representation than the naive attack since an attacker only needs to flip a single bit in the obfuscated intermediate representation, and it is not necessary to figure out any inverse mappings σ_r^{-1} . Also, this attack is independent from the DES structure and can be applied to any round-based block cipher. We try to apply deterministic changes to $\sigma_{n-1}(L_{n-1}, R_{n-1})$, the state going into the last round, and then run the last round operation.

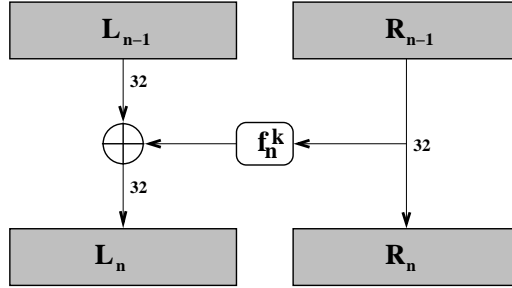


Figure 3: Last round with the round function f_n^k . In the last round the right hand side and the left hand side of the output are usually not crossed over.

Figure 3 shows the last round of the cipher. An attacker knows $R_n = R_{n-1}$ from the ciphertext which is also the input to the round function of the last round. In addition an attacker can modify R_{n-1} even if the mapping of σ_{n-1} is unknown by changing R_n in the ciphertext, decrypting the ciphertext, and encrypting the resulting plaintext afterward. Therefore we have two preconditions for the attack: First, both encryption and decryption operations need to be available, and second, the attacker needs to be able to modify the ciphertext arbitrarily. Using this technique we can find out the positions of $\mu_r(i)$ for $i = 0..32$ which describe the bits for the left-hand side. From the definition of $T_{r,i}^k$, it is clear, that if the attacker keeps the right-hand side input constant, the observed changes in the input to the t-boxes uniquely refer to changes in the left-hand side of the input. The attacker is not able to set L_{n-1} to 0 since she would need to know the round function and hence the key. Therefore, $R_n = 0$ and $L_{n-1} = f_n^k(0) \oplus L_n$.

Now the attacker builds a table of

$$\Delta(c) := \sigma_{n-1}(c, 0) \oplus \sigma_{n-1}(0, 0)$$

for $c = 1..2^{32}$.

Since σ_r contains the unknown non-linear bijection δ_{r-1} it is not possible to build a linear operator in Δ . However, using the table the attacker can always reconstruct the left-hand side of the input in the scenario where the right-hand side is 0. Furthermore, different bits of the left-hand side L_{n-1} can correspond to the same t-box, and in this case the encoding depends on two bits. Therefore, in the first part the attacker tests which bits correspond to the same t-box and then tries all possible bit combinations into this t-box. In this way the attacker gets all possible values for σ_r induced by the left-hand side L_{n-1} . Determining the original value $L_{n-1} \oplus f_n^k(0)$ given the intermediate representation is just a table lookup.

The idea now is to inject faults into the input to the s-box and observe the output. Unfortunately, the attacker does not know how the right-hand side gets encoded in σ_r . In order to get around this

problem the attacker feeds a value x into R_{n-1} that is different from 0 and then resets L_{n-1} to 0. Finally, L_n contains $f_n^k(x) \oplus f_n^k(0)$, and the attacker can extract the key for the last round using differential cryptanalysis. Getting the DES key from the round key requires a 2^8 brute-force search. The problem is that if the right hand side R_{n-1} changes to some value $\neq 0$ the t-box inputs collide with the 16 bits of the left-hand side L_{n-1} . Therefore it is not possible to decode the left-hand side L_{n-1} uniquely since complete new values might show up in the t-boxes that are taking as input bits from the left-hand side.

However, if the attacker sets only one bit in R_{n-1} at most two different t-box outputs are affected, and hence the attacker can simply count the occurrences of the encoded 4-bit values at a certain position in σ_r .

We describe the algorithm for the attack when the specification of the round function is known. We will explain at the end of the algorithm how the algorithm needs to be changed to attack an unknown round function. For convenience we use $D^k(c)$ to describe the decryption of ciphertext c using key k , and $E_i^k(p) = (L_i, R_i)$ to describe iteration of plaintext p for i rounds in the encryption operation using key k . $s^n(k) = s_n^1(k)|\dots|s_n^8(k)$ is the key schedule for key k in round n , m is the size of the input word, n_s the number of s-boxes within the round function, and $sb_n(x) = sb_n^1(x_1)|\dots|sb_n^8(x_8)$:

$$f_n^k(x_1|\dots|x_8) := sb_n^1(x_1 \oplus s_n^1(k))|\dots|sb_n^8(x_8 \oplus s_n^8(k))$$

In our simplified model the in- and outputs of the s-box have the same size, and the system computes the xor of the key and the input to the s-box. The algorithm consists of 3 basic operations: A *Set* operation changes any arbitrary variable. When we do a *Compute* we execute an operation in the iterated block cipher. This can be encryption, decryption, or just a single round of the cipher. *Derive* computes values on known variables without executing the cipher. Figure 4 illustrates the single steps of the algorithm.

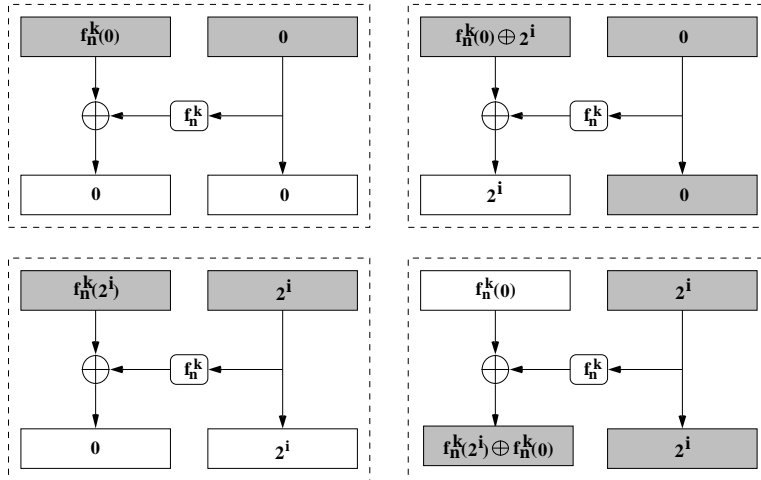


Figure 4: Attacking the last round of the iterated block cipher. Boxes having a white background indicate that the attacker changed values. The picture on the top left shows the initialization of the algorithm (step 1). Afterward, on the top right we change L_n to 2^i in order to reconstruct $\psi_{n-1}(x)$ (step 2). In the bottom left we set 2^i to be input to the round function. The fault injection takes place on the bottom right (step 3): We reset L_{n-1} to $f_n^k(0)$ and obtain the difference $f_n^k(2^i) \oplus f_n^k(0)$ in L_n .

Our attack algorithm works as follows:

1. **Initialization:** (Figure 4 top left)

- (a) Set $L_n := 0, R_n := 0$
 Compute $\sigma_{n-1}(L_{n-1}, R_{n-1}) = E_{n-1}^k(D^k(L_n, R_n))$
 Result: $L_{n-1} = f_n^k(0), R_{n-1} = 0$
 Derive $\Omega = \sigma_{n-1}(L_{n-1}, R_{n-1}) = \sigma_{n-1}(f_n^k(0), 0)$

2. **Reconstruct $\Delta(x)$:** (Figure 4 top right)

- (a) For $i = 1$ to m :
 Set $L_n := 2^i, R_n := 0$
 Compute $\sigma_{n-1}(L_{n-1}, R_{n-1}) = E_{n-1}^k(D^k(L_n, R_n))$
 Set $\Delta(i) = \sigma_{n-1}(L_{n-1}, R_{n-1}) \oplus \Omega$
 For $j = 1$ to $\frac{m}{4}$:
 If $(\Delta(j) \neq 0)$
 Set $\mathcal{O}(j) = \mathcal{O}(j) \cup \{i\}$
- (b) For $j = 1$ to m :
 Set $x = 0$
 For $k = 1$ to $|\mathcal{O}(j)|$:
 Set $e[k] = 0$
 Set $pos = 0$
 For $l = 1$ to $2^{|\mathcal{O}(j)|}$:
 Set $e[pos] = e[pos] + 1 \pmod{2}$
 If $(e[pos] = 0)$
 Set $pos = pos + 1$
 Else
 Set $pos = 0$
 Set $L_n := e[0] \dots e[|\mathcal{O}(j)|], R_n := 0$
 Compute $\sigma_{n-1}(L_{n-1}, R_{n-1}) = E_{n-1}^k(D^k(L_n, R_n))$
 Set $\Delta(L_n) = \sigma_{n-1}(L_{n-1}, R_{n-1}) \oplus \Omega$

3. **Reset L_{n-1} to $f_n^k(0)$:** (Figure 4 bottom left)

- For $i = 1$ to m :
- (a) Set $L_n := 0, R_n := x_i := x_1 \dots x_m$ ($x_i = 1, x_l = 0$ for $l \neq i$)
 Compute $\sigma_{n-1}(L_{n-1}, R_{n-1}) = E_{n-1}^k(D^k(L_n, R_n))$,
 Result: $L_{n-1} = f_n^k(R_n), R_{n-1} = x_i$.
- (b) Derive $w := \sigma_{n-1}(L_{n-1}, R_{n-1}) \oplus \Omega$
 Result: $w = \sigma_{n-1}(f_n^k(R_n), R_n) \oplus \sigma_{n-1}(0, 0)$.
- (c) For x in Δ^{-1}
 For $i = 1$ to 24
 If $\left((\Delta(x)|_{4i}^{4(i+1)} = w|_{4i}^{4(i+1)}) \right)$
 $w := w \oplus \Delta(x)$
- (d) Compute $(L'_n, R'_n) = (\tau_n'' \tau_{n+1}'')(w) = (\sigma_{n-1}^{-1} \pi_n^k)(w)$
 Result: $L'_n \approx f_n^k(x_i) \oplus f_n^k(0), R'_n \approx x_i$

4. **Do differential cryptanalysis to extract the key for the round function f_n^k :**

$l[i] = L'_n |_{4i^{4(i+1)}}$
 $r[i] = EP(R'_n) |_{6i^{6(i+1)}}$
For $s = 1$ **to** n_s :
 (a) **For** $i = 1$ **to** m :
 Compute $c^s[i]$: $sb_n^s(r^s[i] \oplus c^s[i]) = l^s[i]$
 Compute $d^s[i] + +$
 (b) **Set** $c^{s'} := c^s[\max_{i=1}^m d^s[i]]$

5. Reconstruct the original key:

- (a) $k := c^{0'} | c^{1'} | \dots | c^{n_{s'}}$
- (b) **Compute** $s_n(k)^{-1}$ to retrieve original key
- (c) do a brute-force search on the remaining bits of the key.

Step 2 of the algorithm reconstructs $\Delta(x)$, in step 3 we inject the fault by resetting L_{n-1} to $f_n^k(0)$ and computing $L_n = f_n^k(R_n) \oplus f_n^k(0)$. In steps 4 and 5 we compute the key given a round function f_n^k by concatenating the components going into the s-boxes, inverting the key schedule, and running a brute-force search on the remaining key bits.

If the key schedule $s_n(k)$ for round n is unknown, we cannot do step 5 to get the key out. In this case we have to compute the key for round n and then use this key to attack round $n - 1$ until we extract all round keys. If the round function f_i^k is unknown, we can first try out different known round functions (e.g. Skipjack, Blowfish, DES etc) for f_i^k . If none of them works, we have to do cryptanalysis to recover the s-boxes from scratch. We make the basic assumption that the round function is based on an s-box with fixed inputs.

This attack is fully automated and can be run without any knowledge of the system. The attack in steps 1-5 extracts the key in $O(m)$ cryptographic operations, and therefore undermines the security of the obfuscation system.

2.3 Summarizing the attack

We exploit two weaknesses in this attack: First, the boundaries of the rounds are identifiable and protection of intermediate results against tampering is not strong enough. This means that a) hiding the rounds can strengthen the implementation and b) data needs to be safe against leaking of information during execution.

In this attack we show that faults in ciphers are a cheap and efficient technique to extract a secret key from an obfuscated cipher implementation in software. Our attack on obfuscated cipher implementations in software requires only a few cryptographic operations, and therefore an adversary can run the attack on any inexpensive hardware.

We had to modify the original algorithm for differential fault analysis [17] in several steps. The main difference is that it is not possible to inject random faults since the intermediate representation is obfuscated and has multiple points of failure. However, it is still possible to find out a sufficient amount of information about the obfuscated intermediate representation that make it possible for an attacker to inject faults.

In the underlying attack model it is the goal to decrypt some media stream on different machines at the same time. To do this we assume that copy protection of the decryption system is sufficiently strong, and therefore an attacker has to extract the secret key. In the current implementation our attack requires that a decryption system colludes with an encryption system, but actually an attacker

only needs to obtain plaintexts for $2m$ chosen plaintexts and the decryption system. Or, since the system is a symmetric block cipher, we run the attack on the encryption system and need $2m$ chosen ciphertexts from the decryption operation. Furthermore, it is an open question how difficult it is to turn an obfuscated decryption system into an encryption system. In this case having the decryption system is sufficient for the attack.

In the recommended variant the system executes the encryption operation $E'(x) = (f^{-1}Eg)(x)$ and the decryption operation $D'(x) = (g^{-1}Df)(x)$ where f and g are non-linear bijective encodings. The current attack is now impossible, but the disadvantage is that given a ciphertext it is only possible to decrypt when f , g , and the key k are known, or the obfuscated decryption program is being used. It is not implementing DES anymore.

It is crucial to fix the weaknesses in the system or implement other techniques to prevent any common attacks that recover the secret key. In the following sections we explore what we can do about the weaknesses and investigate how to strengthen obfuscation techniques against common attacks.

3 Theoretical Considerations

The weaknesses in this attack are specific to the implementation of the obfuscated cipher. We were able to use specific properties of the DES cipher and the obfuscation method in order to extract the secret key. However, theoretical considerations do not necessarily limit any stronger obfuscation techniques. Here we give a simple argument why the general problem of retrieving embedded data from a circuit is NP-hard, and therefore no efficient general deobfuscator exists for this problem.

In MATCH-FIXED-INPUT we are given two circuits, one of which has additional input k . It is the goal to find a k such that the two circuits are equivalent.

Definition: MATCH-FIXED-INPUT: Given circuits two $C(x, k)$ and $C'(x)$ where $x \in \{0, 1\}^n$ and $k \in \{0, 1\}^c$ where $c \in N$ is constant, find $k' \in \{0, 1\}^c$ such that $\forall x : C(x, k') = C(x)$.

Theorem: MATCH-FIXED-INPUT is NP-hard.

Proof: We reduce SAT to MATCH-FIXED-INPUT which is almost trivial. In order to test satisfiability of circuit $D(x)$, set $C(x, k) = D(k)$ and $C'(x) = true$, and run MATCH-FIXED-INPUT. If MATCH-FIXED-INPUT returns a k' such that $C(x, k') = C'(x)$, then according to the definition there exists an x such that $D(x) = true$. If MATCH-FIXED-INPUT does not return a k' , then for all x $D(x) = false$. Hence, we reduce SAT to MATCH-FIXED-INPUT. \square

For practical purposes, however, this theoretical observation is not much of a relevance since the problem is hard in the worst case but can still be easy for practical purposes. On the average the problem MATCH-FIXED-INPUT is NP-hard, but in several cases heuristic methods can extract the fixed input as in the example of this obfuscated DES cipher.

4 Strengthening Obfuscation

In this section we briefly discuss various mechanisms for defending against our attack using software faults. We first describe some common attacker goals when attacking obfuscated code:

- **Hide data in the program:** The attacker wants to find out certain data values. This case subdivides into the possibility of tracing values during runtime and discovering static values in the code.

- **Protect the program from controlled manipulation:** In this case the attacker wants to force the program to behave in a certain way, e.g. to remove copy protection mechanisms or to cause damage on a system.
- **Hide algorithms of the program:** According to Kerckhoff’s principle cryptographic algorithms are usually public, but in some cases it is useful to hide certain properties by which an attacker can recognize the algorithm, i.e. distinguish for example between AES, IDEA or Blowfish [32, 30, 24].

Often when obfuscating a cipher, commercial tools first encode the plaintext using some hidden encoding function, then run the cipher, and finally decode the ciphertext using some other hidden decoding function. More precisely, the encryption process looks like $E'_k(x) = (F \cdot E_k \cdot G^{-1})(x)$ where E_k is the original DES encryption [21]. Note that F and G must be one-to-one functions so that decryption is possible. The decryption process is similar: $D'_k(x) = (G \cdot D_k \cdot F^{-1})(x)$. This pre- and post-encoding makes chosen ciphertext attacks more difficult since an adversary first needs to recover G . As a result, these encoding makes our fault attack harder to mount. One can still potentially attack the system by using a fault attack against inner levels of the Feistel cipher.

4.1 Defending against a fault-based attack

We mention a few mechanisms for protecting obfuscated systems from a fault attack. One approach is to protect all intermediate results using checksums. These checksums are frequently checked by the obfuscated code. We refer to this approach as *local checking*. Clearly the code for checking these checksums must be hidden in the total program code so that an attacker cannot disable these checkers. One approach for using checksums to ensure code integrity is explained in [15]. In this approach we compute checksums for parts of the program and verify them during program execution. In the extreme we verify a checksum for every single instruction and every data element.

Another approach for checking the computation of obfuscated code is to use *global checking*. The idea is to execute the obfuscated program k times (e.g. $k = 3$) by interleaving the k executions. At the end of the computation the code verifies that all k executions resulted in the same value. As before, the checker must be obfuscated in the code so that it cannot be targeted by the attacker. This global checking approach makes our attack harder since the attacker now has to modify internal data consistently in all k executions of the code.

The problem with the checking approaches is the vulnerability of the checker since it is unprotected against any tampering attack. One approach to make the checker more robust is to obfuscate it and have it verify its own integrity repeatedly while it is checking the program. This variant reduces the maximum time interval an attacker has to run the modified program. In any case the attacker needs to modify to system at more than one place. We note that if the integrity check fails the program should not stop execution immediately since this will tell an attacker where the checker is.

Another approach for making the fault attack more difficult is to diversify the obfuscation mechanism. In other words, each user gets a version of the code that is obfuscated differently (e.g. by using different encoding functions). In diversification we add randomness to the obfuscation methods, and therefore two obfuscated programs are always different after obfuscation. Especially vulnerable places in a program such as the intermediate results of the iterated round-based cipher need to be diversified.

5 Related work

Informally tamper-resistance of a software implementation measures to what extent the implementation resists arbitrary or deliberate modifications. For example, an implementation can be protected

from removing a copy protection mechanism. Thus, obfuscation is a common technique for improving tamper-resistance. Barak et al. [16] give a formal definition of obfuscation using a black-box approach which is the ideal case. They show that in their model, that obfuscation is not possible. Encrypting the executable binary [11] is the most common approach for hiding code. In binary encryption the program is encrypted and decrypts itself during runtime. The problem is that the program is available in the clear at some point before it gets executed on the processor, and it can be intercepted. Furthermore, the system needs to hide the decryption key, and that reduces recursively to the key obfuscation problem itself.

A common approach for obfuscation is to obstruct common static program analysis [35, 22, 34]. The main technique for doing this is to insert of additional code that creates pointer aliasing situations. Applying static program analysis to analyze a program containing possible pointer aliasing turns out to be NP-hard [29]. This obfuscation technique only protects against attacks by static program analysis. It is still possible to do dynamic attacks with a debugger or any type of tampering.

The goal of obfuscation is to hide as many program properties as possible. The principle of improving tamper-resistance by obfuscation is that if an attacker cannot find the location for manipulating a value, it is impossible to change this value. In addition an obfuscator can eliminate single points of failure. On the other hand obfuscation never protects against existential modification.

Collberg et al define some metrics for obfuscation in [22]. They classify obfuscation schemes by the confusion of a human reader (“potency”), the successfulness of automatic deobfuscation (“resilience”), the time/space overhead (“cost”), and the blending of obfuscated code with original code (“stealth”). But obfuscation of a secret key requires stronger properties of obfuscation, since any definition of tamper-resistance is missing. A program that is a good obfuscator in these metrics can still have a single point of failure, and therefore it does not protect the program against fault attacks.

Tamper-resistance can also be improved by techniques other than obfuscation. We already mentioned self-checking of code as one possibility [15, 27, 9]. Protection by software guards is another technique to prevent tampering [20]. Software guards are security modules that implement different tasks of a program and thus eliminate single points of failure. In addition a program can implement anti-debugging techniques in order to prevent tampering with a debugger [19]. Anti-debugging inserts instructions into a program or changes properties in order to confuse a debugger. For example a program can arbitrarily set break points or misalign code. Furthermore, virtual software processors are a technique for making tampering difficult [12]. Virtual software processors run the original program on a software processor, and in order to reverse engineer the original program, an attacker needs to compromise any protection mechanism of the virtual software processor as well.

Goldreich and Ostrovsky show in [26] that software protection against eavesdropping can be reduced to *oblivious simulation of RAMs*. In their definition a RAM is oblivious if two different inputs with the same running time create equivalent sequences of memory accesses. Oblivious RAM protects against any passive attack and therefore strengthens an obfuscator because it is impossible to find out the memory locations a program accesses. However, it does not protect against the fault injection attack.

Current hardware dongles are based on the idea of oblivious RAM, since the code implementing the license check sits on the dongle.

6 Open Problems

In other areas of information hiding techniques, such as watermarking, benchmark programs are available to measure the strength of a technique to hide information. For example, StirMarks [31] uses a variety of different generic attacks on a watermarked image to make the watermark illegible. It is an open problem to build such a benchmark for code obfuscation and tamper resistance tools. Such a benchmark would take as input some tamper resistant code and attempt to break the tamper resistance. Currently no such benchmark exists and there is no clear model for building such a benchmark.

One of the main open problems in code obfuscation is to come up with a model for obfuscation that can be realized in practice. [16] defines obfuscation using a black-box model that hides all properties of a program. They show that it is not possible to achieve obfuscation in that model. For practical purposes a black box model might not always be necessary. In the example of the obfuscated DES cipher in this paper we only need to make sure that it is impossible to get information about the secret key. The open research problem is to find the most general definition for obfuscation that can be realized in practice.

7 Conclusion

Code obfuscation provides some protection against attackers who want to find out secret data or properties of a program, but it is not sufficient as a stand-alone system. In this study we evaluate the usability of obfuscation when hiding a secret key in an iterated round-based software cipher. We find weaknesses in a commercial state-of-the-art obfuscator. Our attack enables automated extraction of the secret key from the obfuscated program code. We discuss a few methods for defending against these attacks.

References

- [1] Adobe EBooks. <http://www.adobe.com/epaper/ebooks>.
- [2] CSS. <http://www.dvdcca.org/css>.
- [3] Intertrust. <http://www.intertrust.com>.
- [4] Microsoft Windows Media Technologies. <http://www.microsoft.com/windows/windowsmedia>.
- [5] RetroGuard Java Obfuscator. <http://www.retrologic.com>.
- [6] TCPA. <http://www.trustedpc.org>.
- [7] Soft microcontroller data book, 1993. Dallas Semiconductor.
- [8] Cloakware Corporation, World Intellectual Property Organization, WO 00/77596 A1, 2000.
- [9] Intel Corporation, US Patent Office, US 6,205,550, 2000.
- [10] Intertrust Corporation, US Patent Office, US 6,157,721, 2000.
- [11] Armouring the ELF: Binary encryption on the UNIX platform. *Phrack Inc.*, (58), 2001.
- [12] Microsoft Corporation, World Intellectual Property Organization, WO 02/01327 A2, 2002.

- [13] D. G. Abraham, G. M. Dolan, G. P. Double, and J. V. Stevens. Transaction Security System. *IBM Systems Journal*, 30(2):206–229, 1991.
- [14] R. Anderson and M. Kuhn. Low cost attacks on tamper resistant devices. In *Proc. 5th International Security Protocols Conference*, pages 125–136, 1997.
- [15] D. Aucsmith. Tamper-resistant software: An implementation. *Lecture Notes in Computer Science*, 1174:317–333, 1996.
- [16] B. Barak, O. Goldreich, R. Impagliazzo, S. Rudich, A. Sahai, S. Vadhan, and K. Yang. On the (im)possibility of obfuscating programs. *Lecture Notes in Computer Science*, 2139:1–18, 2001.
- [17] E. Biham and A. Shamir. Differential fault analysis of secret key cryptosystems. *Lecture Notes in Computer Science*, 1294:513–525, 1997.
- [18] D. Boneh, R. A. DeMillo, and R. J. Lipton. On the importance of checking cryptographic protocols for faults. *Lecture Notes in Computer Science*, 1233:37–51, 1997.
- [19] S. Cesare. Linux anti-debugging techniques. *Security Focus*, Jan. 1999.
- [20] H. Chang and M. J. Atallah. Protecting software code by guards. In *Proc. of Workshop on Security and Privacy in Digital Rights Management 2001*. Association of Computing Machinery.
- [21] S. Chow, H. Johnson, P. C. van Oorschot, and P. Eisen. A White-Box DES Implementation for DRM Applications. In *ACM CCS-9 Workshop DRM 2002*.
- [22] C. Collberg, C. Thomborson, and D. Low. Manufacturing cheap, resilient, and stealthy opaque constructs. In *The 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '98)*, pages 184–196, New York, Jan. 1998. Association for Computing Machinery.
- [23] S. A. Craver, M. Wu, B. Liu, A. Stubblefield, B. Swartzlander, D. S. Wallach, D. Dean, and E. W. Felten. Reading between the lines: Lessons from the SDMI challenge. In *Proc. 10th USENIX Security Symp.*, 13–17 Aug. 2001.
- [24] J. Daemen and V. Rijmen. Rijndael for AES. In NIST, editor, *The Third Advanced Encryption Standard Candidate Conference, April 13–14, 2000, New York, NY, USA*, pages 343–347, Gaithersburg, MD, USA, 2000. National Institute for Standards and Technology.
- [25] W. Diffie and M. Hellman. New directions in cryptography. *IEEE Transactions on Information Theory*, IT-22(6):644–654, Nov. 1976.
- [26] O. Goldreich and R. Ostrovsky. Software protection and simulation on oblivious RAMs. *Journal of the Association for Computing Machinery*, 43(3):431–473, May 1996.
- [27] B. Horne, L. Matheson, C. Sheehan, and R. E. Tarjan. Dynamic self-checking techniques for improved tamper-resistance. In *Proc. of Workshop on Security and Privacy in Digital Rights Management 2001*. Association of Computing Machinery.
- [28] P. Kocher, J. Jaffe, and B. Jun. Differential power analysis. *Lecture Notes in Computer Science*, 1666:388–397, 1999.
- [29] W. Landi. Undecidability of static analysis. *ACM Letters on Programming Languages and Systems*, 1(4):323–337, Dec. 1992.

- [30] A. J. Menezes, P. C. Van Oorschot, and S. A. Vanstone. *Handbook of applied cryptography*. CRC Press, 1997.
- [31] F. A. P. Petitcolas, R. J. Anderson, and M. G. Kuhn. Attacks on copyright marking systems. *Lecture Notes in Computer Science*, 1525:219–239, 1998.
- [32] B. Schneier. *Applied Cryptography*. Wiley, 1994.
- [33] A. Shamir and N. van Someren. Playing “hide and seek” with stored keys. *Lecture Notes in Computer Science*, 1648:118–124, 1999.
- [34] B. Steensgaard. Points-to analysis in almost linear time. In *POPL’96*, pages 32–41. ACM Press, Jan. 1996.
- [35] C. Wang, J. Davidson, J. Hill, and J. Knight. Protection of software-based survivability mechanisms. Proceedings of the 2001 Dependable Systems and Networks (DSN’01), 2001.