

# An Information Theoretic Perspective on the Differential Fault Analysis against AES

Yang Li, Shigeto Gomisawa, Kazuo Sakiyama, Kazuo Ohta

The University of Electro-Communications  
1-5-1 Chofugaoka, Chofu-shi, Tokyo 182-8585, Japan  
{liyong, g-shigeto-lfat, saki, ota}@ice.uec.ac.jp

**Abstract.** Differential Fault Analysis against AES has been actively studied these years. Based on similar assumptions of the fault injection, different DFA attacks against AES have been proposed. However, it is difficult to understand how different attack results are obtained for the same fault injection. It is also difficult to understand the relationship between similar assumptions of fault injection and the corresponding attack results. This paper reviews the previous DFA attacks against AES based on the information theory, and gives a general and easy understanding of DFA attacks against AES.

We apply the similar analysis on DFA attacks on AES-192 and AES-256, and we propose the attack procedures to reach the theoretical minimal number of fault injections.

**Keywords:** Differential Fault Analysis, AES, Information theory

## 1 Introduction

In 1997, public key cryptosystems were pointed to be vulnerable to fault attacks that use the computational errors during the execution to find the secret key [7]. At the same year, Biham and Shamir applied this idea to block cipher DES and introduced the concept of Differential Fault Analysis (DFA) [5]. Given an encryption of a block cipher, a fault-free ciphertext can be obtained for a plaintext. Then by injecting a certain kind of fault during the execution of the cryptographic calculation, attackers can get a faulty ciphertext as well. The assumptions of the fault injection are referred as the **fault model** in this paper. DFA obtains the information of the secret based on the fault-free ciphertext and the faulty ciphertext under a certain fault model.

This paper focuses on the DFA attacks against Advanced Encryption Standard (AES) excluding the cases where faults are injected at the key schedule. In the early stage of DFA attacks against AES, the full recovery of the secret key were likely to require 50 to 250 faulty ciphertexts [8, 9, 13]. Later in researches shown in [10, 12, 14, 15], only two or one faulty ciphertext was enough for the full recovery of the secret key. However, these papers presented different attack results for the same fault model that disturbs a single byte. In 2009, a DFA attack called diagonal fault analysis was proposed [6]. Their attack can retrieve the full key with one

faulty ciphertext with a fault model allowing multiple faulty bytes. In 2006, Moradi, Shalmani and Salmasizadeh proposed a generalized DFA attack against AES [11]. In brief, the generality of their attack is achieved by dividing all possible faults into two groups, and giving attacks for each group. By simulation, they found that for the first group, 6 faulty ciphertexts in average can identify the secret key, while 1500 faulty ciphertexts are needed for the second group. It is not difficult to find that the fault models used in [6] and [11] are also similar to that used in [10, 12, 14, 15], however different attack results are obtained as well.

In this paper, DFA attacks against AES are analyzed through the observation of information theory. Assumptions in a fault model are regarded as the information of the difference between a fault-free intermediate value and a faulty intermediate value. Based on the relationship between the difference of intermediate values and the secret key, we give a simple understanding for the existing DFA attacks against AES. Our analyses find that there is a limitation of the attack efficiency for each fault model. Attacks in [10], [15], [6] and [11] reached the limitations of their fault models, while attacks in [12, 14] did not. Also, we propose a simple model for predicting the attack efficiency. Our prediction model obtains the similar results with the simulation results provided by previous papers. In a word, this paper provides a generalized and simple understanding for the DFA attacks based on the similar fault models, and proposes an optimized attack flow for DFA attacks as well.

For the DFA attacks on AES-192 and AES-256, we manage to predict the minimal times of the fault injection based on the information theory. Then we also propose the DFA attack procedures that can use the minimal times of fault injections to retrieval the key within a practical computational complexity. The proposed DFA attack on AES-192 has been verified by the simulations based on C language.

This paper is organized as follows. In Sect. 2, we briefly explain the structure of AES. In Sect. 3, we generally analyze the DFA attacks through the perspective of information theory. In Sect. 4, several related previous DFA attacks against AES-128 are explained based on the observation of information theory. In Sect. 5, DFA attacks on AES-192 and AES-256 are analyzed and discussed. In Sect. 6, we discuss the possible future research about the DFA attacks against AES. In Sect. 7, we conclude this paper.

## **2 Overview of the Structure of AES-128**

In 2000, AES was selected as the new standard of symmetric key encryption scheme by the US government. AES is 128-bit block cipher, and has three kinds of key sizes as 128, 192 and 256 bits. Except Sec. 5, the discussion in this paper is based on the

encryption of AES with 128-bit secret key. In this paper a plaintext, an intermediate value and a ciphertext of AES-128 are denoted by  $P$ ,  $I$  and  $C$ , respectively. AES operates on a  $4 \times 4$  state matrix as shown in Table 1. Every element of the state matrix is a byte represented by  $I_{ij}$ , where  $i, j \in [0, 3]$  and  $i, j$  are its row and column positions, respectively. Notice that,  $P$ ,  $K$  and  $C$  can be expressed in the same manner.

$I_{00}$	$I_{01}$	$I_{02}$	$I_{03}$
$I_{10}$	$I_{11}$	$I_{12}$	$I_{13}$
$I_{20}$	$I_{21}$	$I_{22}$	$I_{23}$
$I_{30}$	$I_{31}$	$I_{32}$	$I_{33}$

**Table 1.** AES state matrix

AES-128 consists of 10 rounds. Each round has its own round key denoted by  $K^i$ , where  $i \in [0, 10]$ . Each round key can be expanded from the original key  $K$  by the AES key schedule scheme. Conversely, obtaining a round key is equivalent to obtaining the original key and all the other round keys. After the initial AddRoundKey, the first 9 rounds of AES are the same consisting of four AES operations as SubBytes, ShiftRows, MixColumns and AddRoundKey. The last round of AES-128 only consists of SubBytes, ShiftRows and AddRoundKey.

Before we introduce the details, we list several notations used in this paper. We denote the faulty intermediate value and the faulty ciphertext by  $I'$  and  $C'$ , respectively. The difference between  $I$  and  $I'$  is denoted by  $\Delta I$ , and that between  $C$  and  $C'$  is denoted by  $\Delta C$ .

The functionalities of AES operations on the real values (*e.g.*  $I$ ,  $I'$ ) and differences (*e.g.*  $\Delta I$ ) are briefly explained as follows.

#### AddRoundKey(ARK)

As shown in Table 2, the AddRoundKey performs the exclusive OR calculation ( $\oplus$ ) between the current state and the corresponding round key. AddRoundKey affects the real value of each byte in the state, but does not affect the difference between the fault-free state and faulty state.

#### SubBytes(SB)

Each byte of the state is substituted by another value according to the pre-computed S-box table. The mapping of AES S-box is bijective.

#### ShiftRows(SR)

As shown in Table 3, the rows of the state are cyclically shifted according to the row number, so does the differences.

$I_{00}$	$I_{01}$	$I_{02}$	$I_{03}$	$I_{00} \oplus K_{00}$	$I_{01} \oplus K_{01}$	$I_{02} \oplus K_{02}$	$I_{03} \oplus K_{03}$
$I_{10}$	$I_{11}$	$I_{12}$	$I_{13}$	$I_{10} \oplus K_{10}$	$I_{11} \oplus K_{11}$	$I_{12} \oplus K_{12}$	$I_{13} \oplus K_{13}$
$I_{20}$	$I_{21}$	$I_{22}$	$I_{23}$	$I_{20} \oplus K_{20}$	$I_{21} \oplus K_{21}$	$I_{22} \oplus K_{22}$	$I_{23} \oplus K_{23}$
$I_{30}$	$I_{31}$	$I_{32}$	$I_{33}$	$I_{30} \oplus K_{30}$	$I_{31} \oplus K_{31}$	$I_{32} \oplus K_{32}$	$I_{33} \oplus K_{33}$

**Table 2.** AES AddRoundKey

$I_{00}$	$I_{01}$	$I_{02}$	$I_{03}$	$\xrightarrow{SR}$	$I_{00}$	$I_{01}$	$I_{02}$	$I_{03}$
$I_{10}$	$I_{11}$	$I_{12}$	$I_{13}$		$I_{11}$	$I_{12}$	$I_{13}$	$I_{10}$
$I_{20}$	$I_{21}$	$I_{22}$	$I_{23}$		$I_{22}$	$I_{23}$	$I_{20}$	$I_{21}$
$I_{30}$	$I_{31}$	$I_{32}$	$I_{33}$		$I_{33}$	$I_{30}$	$I_{31}$	$I_{32}$

**Table 3.** AES ShiftRows

### MixColumns(MC)

A linear transformation performed on each column of the state computed by

$$\begin{bmatrix} I_{0j} \\ I_{1j} \\ I_{2j} \\ I_{3j} \end{bmatrix} = \begin{bmatrix} 02 & 03 & 01 & 01 \\ 01 & 02 & 03 & 01 \\ 01 & 01 & 02 & 03 \\ 03 & 01 & 01 & 02 \end{bmatrix} \cdot \begin{bmatrix} I_{0j} \\ I_{1j} \\ I_{2j} \\ I_{3j} \end{bmatrix},$$

where  $j \in [0, 3]$  and the multiplication is performed in  $\text{GF}(2^8)$ . Since MixColumns is a linear transformation, we have  $MC(I) \oplus MC(I') = MC(I \oplus I') = MC(\Delta I)$ .

## 3 Information Theoretic Perspective on DFA

### 3.1 The Fault Model for DFA

Every fault model assumes that it is possible to inject a certain type of faults at a certain state of the AES calculation. For example, the most discussed fault model in this paper is the one that assumes attackers can disturb a random byte at the input of the 8<sup>th</sup> round of AES. Hereafter we refer this fault model as **Piret's Fault model**, since it was first introduced in Piret *et al.*'s paper [14]. And we refer the state where a fault is injected by the **injection state**. Through the perspective of information theory, the fault model can be considered as the information of  $\Delta I$  at the injection state (the difference between the fault-free intermediate value and the faulty one). In addition to the values of  $(C, C')$ , DFA is the cryptanalysis that obtains the information of the secret key based on the information of  $\Delta I$  at the injection state.

The difference between two different random 128-bit values have  $2^{128} - 1$  candidates. While in Piret's fault model, the difference only have  $255 \times 16$  candidates, where 255 and 16 correspond to 255 possible difference values and 16 possible fault positions. We can see that this fault model provides  $-\log_2 \frac{1}{2^{128}} - (-\log_2 \frac{1}{255 \times 16}) = 116$  bits of information of the  $\Delta I$  at the injection state. In Piret *et al.*'s paper, based on a pair of  $(C, C')$ , the key space of AES-128 can be restricted to  $2^{40}$ . This result can be understood as the attacker obtains  $128 - 40 = 88$  bits information of the secret key from a pair of  $(C, C')$  and  $128 - 12 = 116$  bits information of  $\Delta I$  at the injection state.

### 3.2 Basic Attacks for DFA

In the sense of information theory, given a pair of plaintext and ciphertext  $(P, C)$ , the secret key  $K$  can be identified theoretically. As far as our knowledge, for the full-round of AES-128, the exhaustive search is the most effective way to reveal the key information based on  $(P, C)$ . In the exhaustive search, every possible key candidate is used to encrypt the plaintext to get a ciphertext candidate. Only when the tested key is the correct one, the obtained ciphertext candidate matches the given ciphertext. However, the exhaustive search over  $2^{128}$  key candidates cannot be performed in a practical time.

The exhaustive search for DFA based on  $(C, C')$  and the information of  $\Delta I$  at the injection state also exists. The similar idea was explained in Piret *et al.*'s paper as the **basic attack** for DFA [14]. Under a certain fault model, attackers need to get the correct ciphertext  $C$  and the faulty ciphertext  $C'$ . Every pair of  $(C, C')$  and the fault model can provide information of the key and restrict the key space. Repeat restricting the key space based on new faulty ciphertexts, finally the key can be identified. The algorithm of the basic attack for DFA is shown as follows.

1. Have a guess of the secret key  $K_g$  from a list of possible keys.
2. Calculate the values of  $I$  and  $I'$  at the fault injection state based on  $(C, K_g)$  and  $(C', K_g)$ , respectively.
3. Calculate the difference  $\Delta I = I \oplus I'$  and check whether  $\Delta I$  satisfies the fault model or not. If not, delete  $K_g$  from the key list. Otherwise, keep it in the list.
4. If the key list has more than one candidate, take another faulty ciphertext and repeat steps 1, 2 and 3 to restrict the current key list. Hereafter we refer steps 1, 2 and 3 as a **DFA search**.

On the one hand, for a large key space such as AES-128, this basic attack algorithm is not practical with regard to the computational cost. On the other hand, the

basic attack can fully use the information provided by injected faults in any fault model to restrict the key space, so that it reaches the max attack efficiency with regard to the information theory.

### 3.3 Divide and Conquer used in DFA Attacks

The basic technique used for turning the basic attack into a practical DFA attack is **divide and conquer**. By dividing the 128-bit key into several parts and analyze them part by part, the key search space can be reduced dramatically.

The last three rounds of AES with disturbing one byte at the beginning of the 8<sup>th</sup> round are shown in Fig. 1.

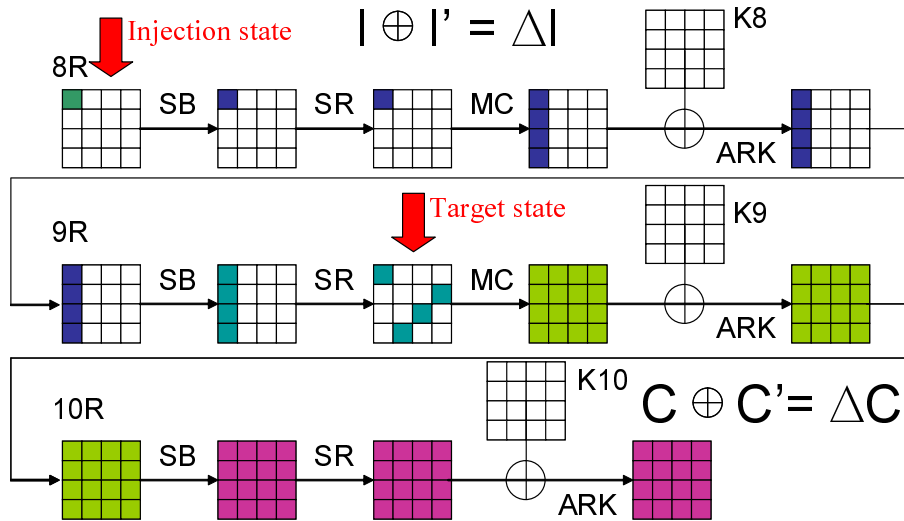


Fig. 1. The last three rounds of 128-bit AES with disturbing one byte at the beginning of the 8<sup>th</sup> round.

One byte fault injected at the beginning of the 8<sup>th</sup> round propagates to four bytes of its column in the 8<sup>th</sup> MixColumns. And these four non-zero faulty bytes will be reserved after the 9<sup>th</sup> SubBytes, and shifts to each column in the 9<sup>th</sup> ShiftRows. Notice that from the 9<sup>th</sup> MixColumns to the output of ciphertext, only the 9<sup>th</sup> MixColumns is the operation relates to 4 bytes, other operations are all byte-wise independent. Here we refer the state before the 9<sup>th</sup> MixColumns as **the target state**. With a guess of four bytes of  $K^{10}$ , four bytes of differential at the target state can be calculated. As a result, the DFA search of 16 bytes of  $K^{10}$  can be divided into

four groups and be searched separately. Each group has four bytes, so that the total search space becomes  $2^{32} \times 4 = 2^{34}$ .

Under the same fault model, as long as different amounts of information from the fault model are used, different attack results will be obtained. The more information of  $\Delta I$  at the injection state used in the DFA search, the less key will be left, as a result the attack will be more efficient in the sense of information theory. However, not all the information of  $\Delta I$  can be easily exploited considering the computational cost. There is a limitation for improving the attack efficiency of DFA in each fault model, which can be achieved by the basic attack theoretically but not practically.

## 4 Review the Previous DFA attacks against AES-128

In 2003, Piret and Quisquater proposed a DFA attack against AES based on disturbing one random byte of the AES state between the 7<sup>th</sup> MixColumns and 8<sup>th</sup> MixColumns [14]. We refer this attack by Piret's attack in this paper. According to their analysis method, two well-located faults are needed for easy retrieving of the key, and one fault can reduce the size of the key space to  $2^{40}$ . According to the structure of AES shown in Fig. 1, one non-zero byte of  $\Delta I$  at the injection state will propagate to 4 non-zero bytes of  $\Delta I$  at the target state. Moreover, each column of  $\Delta I$  at the target state will have one non-zero byte. Then 16 bytes of  $K^{10}$  are divided into four groups to perform the DFA search independently and DFA search checks whether a column of  $\Delta I$  at the target state has only one non-zero byte. The simulations shows that the average size of the key candidates of 4 key bytes is  $2^{10}$ , so that the total number of the key candidates is about  $2^{40}$ . Note that some of the information of  $\Delta I$  at the target state, such as the positions of 4 non-zero bytes are not used in Piret's attack.

In 2009, Mukhopadhyay proposed a DFA similar to Piret's attack [12]. Under the Piret's fault model, it is considered that attackers can first guess where the faulty byte is injected at the injection state. Then the positions of faulty bytes in  $\Delta I$  at the target state are fixed, so that the total key space can be restricted to  $2^{32}$ . Since there are 16 possibilities of the original faulty byte position, the total key space for Piret's fault model can be restricted to  $2^{32} \times 16 = 2^{36}$ . The improvement comes from using the information about the positions of the propagated faulty bytes at the target state. However, this work still does not fully exploit the information of Piret's fault model.

Later, in the same year, Tunstall and Mukhopadhyay further improved the DFA attacks based on the Piret's fault model [15]. At the first step, they guess the fault position and get the key space with size  $2^{32}$ . Then in the second step, they applied

the key schedule scheme to obtain  $K^9$  based on each key candidate of  $K^{10}$ . Then each key candidate can be checked whether it comes from the one faulty byte before the 8<sup>th</sup> MixColumns. For each position, the key space can be restricted to  $2^8$ , so that the total key space can be reduced to  $2^{12}$  in the second step. We can see that this work uses all the information of Piret's fault model and reaches the limitation of Piret's fault model.

In 2010, Gomisawa *et al.* proposed very similar DFA attacks based on the Piret's fault model [10]. First, they found that the positions of 4 faulty bytes at the target state have only 4 patterns instead of 16 patterns analyzed in [12]. So instead of guessing of the exact fault injection position, attackers can guess which position pattern the injected fault belongs to. As a result, the key space can be first restricted from  $2^{128}$  to  $2^{34}$  in the first step. Then by using the information of the fault values of  $\Delta I$  at the target state, the key space can be restricted to  $2^{12}$ . In average, the search algorithm proposed by Gomisawa *et al.* is faster than the algorithm proposed in [15]. For example, the worst case of Gomisawa's algorithm performs searches over  $2^{32}$  keys for 4 times, while that of Tunstall's algorithm is 16 times. For details, we refer to [10].

In 2009, Saha, Mukhopadhyay and RoyChowdhury proposed a Diagonal Fault Attack [6]. Their fault model is that multiple faults are injected at the diagonal of the state matrix at the beginning of the 8<sup>th</sup> round. In their analysis, when only one diagonal is with fault, a pair of  $(C, C')$  can restrict the total key space to  $2^{34}$ . We can see that this result is the same with the intermediate result in [10] and it is already the limitation for this attack model. Different from Piret's fault model, this fault model cannot provide any information about the fault values at the target state.

In 2006, Moradi *et al.* proposed a generalized method of Differential Fault Analysis against AES [11]. In their analysis, all possible faults are divided into two groups. Take a column of  $\Delta I$  at the target state, if at least one of the 4 bytes are fault-free, then this fault belongs to the first group, otherwise, it belongs to the second group. For the first group, the corresponding 32 bits of secret key can be obtained by a fault-free ciphertext and 6 faulty ciphertexts, while it need approximately 1500 faulty ciphertexts to identify 32 key bits for the second group. The fault model of their paper only provides information of the number of faulty bytes in a column of  $\Delta I$  at the target state. The information of their fault model is already fully exploited in their analysis.

The attack results of these DFA attacks against AES are summarized in Table 4.



	Fault type	Attack efficiency in average
[14]	1-byte fault before 8 <sup>th</sup> MixColumns	2 <sup>40</sup> key candidates for a pair of (C, C')
[12]	1-byte fault before 8 <sup>th</sup> MixColumns	2 <sup>36</sup> key candidates for a pair of (C, C')
[10]	1-byte fault before 8 <sup>th</sup> MixColumns	2 <sup>12</sup> key candidates for a pair of (C, C')
[15]	1-byte fault before 8 <sup>th</sup> MixColumns	2 <sup>12</sup> key candidates for a pair of (C, C')
[6]	1 random faulty diagonal before 8 <sup>th</sup> SubBytes	2 <sup>34</sup> key candidates for a pair of (C, C')
[11]	Non-full active column of $\Delta I$ at target state	1 C and 6 C' to identify a 128-bit key
	full-active column of $\Delta I$ at target state	1 C and 1500 C' to identify a 128-bit key

**Table 4.** The summary of attack results of DFA attacks against AES-128.

#### 4.1 The General Attack Flow of DFA

The attack flow of DFA can be mainly divided into two kinds based on whether plaintext is used. If the plaintext corresponding to the fault-free ciphertext is unknown, only faulty ciphertexts can be used to identify the key. Otherwise, attackers can first restrict the key space to a reasonable size based on faulty ciphertexts, and then apply the exhaustive search based on  $(P, C)$  to identify the correct key. We can express this attack flow of DFA as follows.

$$2^{128} \xrightarrow{C, C', \Delta I} 2^{??} \xrightarrow{P, C} 1.$$

The better DFA attacks should request fewer faulty ciphertexts and cost less computations. As a result, better DFA attacks should use more information of every pair of  $(C, C')$  and have a reasonable computational cost at the same time. When a fault model is given, attackers directly obtain the information of  $\Delta I$  at the injection state. Then attackers need to convert the information at the injection state to the one at the target state. Different types of information at the target state cost differently in DFA searches. We try to propose the best attack flow of DFA making a good trade-off between them.

First we separate the information that can be used in a DFA attack into four types as follows.

1. The number of non-zero bytes in each column of  $\Delta I$  at the target state.
2. The positions of non-zero bytes of  $\Delta I$  at the target state.
3. The relationship between values of non-zero bytes of  $\Delta I$  at the target state.
4. The information of  $(P, C)$ .

The first type of information can be exploited by applying divide and conquer and it is the most important information that makes DFA possible. The second type of information can be exploited by arranging the attack results after exploiting the

first type of information. Then since checking the third type of information needs to pass at least two MixColumns, two SubBytes and key schedule, so that divide and conquer cannot be easily applied. The last information can identify the key, but it is the most costly calculation. When these four types of information are all available to attackers, the best attack flow of DFA should first use the first two types of information to restrict key space to a reasonable size. Then, the third type of information can be applied to further restrict the key space. Finally, the last information can be used to identify the key.

## 4.2 Predicting The Attack Efficiency of DFA

In this section, we discuss the relationship between the information of each fault model and the information of  $K$ . According to the structure of AES,  $(I, I')$  at the target state goes through MixColumns, SubBytes, ShiftRows and AddRoundKey to become  $(C, C')$ . Since MixColumns and SubBytes are bijective mapping with regard to a column of state or the entire state, and ShiftRows only change the positions of faulty bytes, we simplify this transformation as

$$BM(I) \oplus K = C, \quad (1)$$

$$BM(I \oplus \Delta I) \oplus K = C', \quad (2)$$

where  $BM$  stands for a bijective mapping, and  $I$  can be a column of target state or the entire target state.

Based on Eq. (1), when  $C$  is fixed, for each value of  $I$ , there is a corresponding value of  $K$ . The key space after a DFA search is equivalent to the number of  $I$  that can pass the Eq (3), where  $\Delta C$  is fixed by  $(C, C')$  and  $\Delta I$  at the target state are restricted by the information from the fault model.

$$BM(I) \oplus BM(I \oplus \Delta I) = \Delta C, \quad (3)$$

For each possible value of  $\Delta C$ , the space of  $\Delta I$  has been restricted by the differential distribution table of  $BM$ , but the key space has not been restricted. After that, when we use the information from the fault model to further restrict the space of  $\Delta I$  at the target state, the key space begins to be restricted.

Base on two conditions, we get a conclusion that the information of  $\Delta I$  at the target state provides the same amount of information to the key. First, we assume that the information of  $\Delta I$  at the target state provided by fault model is independent from that provided by the differential distribution table of  $BM$ . Then, assume the restriction condition of a DFA search covers  $q\%$  of all possible faults, each possible

$\Delta I$  that passed the differential distribution table of  $BM$  have the same probability of  $q\%$  to pass the restriction of the fault model. Second, we assume that the value of  $I$  are uniformly distributed to the values of possible  $\Delta I$  that passed the differential distribution table of  $BM$ . As a result,  $q\%$  of  $I$  will pass the restriction of both the differential distribution table of  $BM$  and the fault model. Finally,  $q\%$  of  $K$  are left after the DFA search.

In the case of the introduced DFA attacks against AES, the used fault model should have little correlation between the differential distribution table of  $BM$ . And according to the differential distribution table of AES S-box, the values of  $I$  are almost uniformly distributed for each possible  $\Delta I$  as well. In a relaxed environment, we can use the conclusion that the size of  $\Delta I$  restricted by the conditions for a DFA search is the same with the size of the key space after this DFA search. In other words, DFA attacks against AES can obtain the same amount of information about  $K$  with the information about  $\Delta I$  used in this attack. It is checked that this prediction matches the simulation results given in [6, 10–12, 14, 15].

For example, we try to predict the attack efficiency of the second type of DFA attacks in [11]. Since the fault model only has the information that four bytes of a column of  $\Delta I$  at the target state are all non-zero, a pair of  $(C, C')$  provides  $\log_2(\frac{2^{32}-1}{255^4}) \simeq 0.02259$  bit information of 4 bytes of  $\Delta I$  at the target state. According to our prediction, the key also obtains about 0.02259 bit information from a pair of  $(C, C')$ . As a result, at least 1420 faulty ciphertexts are needed to recover 32 bits of key ( $1420 \times 0.02259 \simeq 32$ ), where the simulation result in [11] show that in average 1500 ciphertexts can identify 32-bit key. The simulation result also indicates that the information provided by different faulty ciphertexts are almost independent from each other.

## 5 Observations of DFA attacks on AES-192 and AES-256

The DFA attacks on AES-192 and AES-256 were previously discussed in [1–3] and [4]. The latest and best previous attack results against AES-192 and AES-256 are the ones discussed in [4]. In [4], the AES-192 key can be identified with two pairs of correct and faulty ciphertexts and a pair of plaintext and ciphertext. While the AES-256 key can be identified with three pairs of correct and faulty ciphertexts and a pair of plaintext and ciphertext.

Based on the previous analysis based on information theory, we can conclude that the 1-byte random fault with unknown position can provide 116-bit information of the secret key. And the 1-byte random fault with known position can provide 120-bit information of the secret key.

Intuitively, when considering the DFA attacks on AES-192 and AES-256 where the secret information is 192-bit and 256-bit, respectively, the minimal number of required fault injection is two. With two fault injections with known position, the maximal obtainable secret information is  $120 \times 2 = 240$  bits. For AES-192, the secret key should be able to be identified without further exhaustive search. For AES-256, two fault injections should be able to restrict the key space to  $2^{16}$  that is small enough for the exhaustive search. Compared to the attack results in [4], for AES-192, attackers are not required to know the plaintext anymore. For AES-256, attackers can reduce 1 pair of correct and faulty ciphertexts for recovering the key.

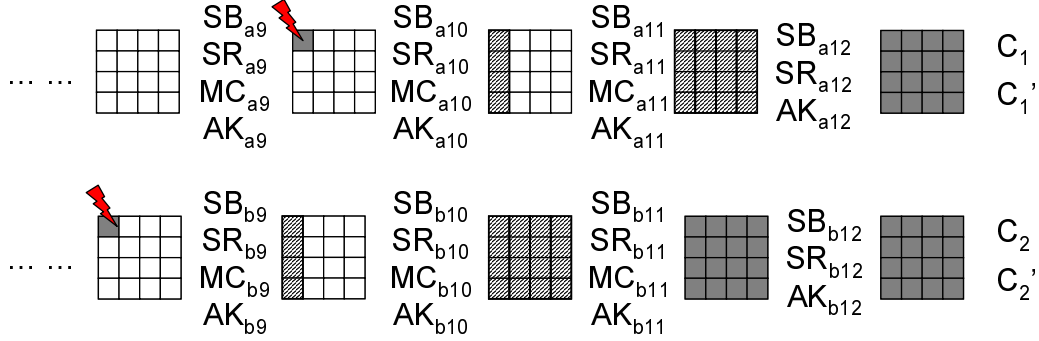
For the simplicity, we show the practical calculation procedures to identify the key for AES-192 and AES-256 with fault injection with known position. Some techniques well used in rebound attacks [16] are used to speed up the key retrieval process. Our attack results are summarized in Table 5. All the introduced attack procedures can be extended for the fault injections with unknown positions. The only penalty is an increase of the calculation complexity. Note that the attack procedures introduced in Sec. 5.1 and Sec. 5.3 are based on the ones introduced in [4]. The difference is that the final attack procedures are checking the deeper differentials in our attacks instead of the exhaustive search in [4].

**Table 5.** Summary of our DFA attack results

AES-192			
Fault Model	No. of Fault	Exhaustive search size	Reference
1 byte random fault	2	1	Sec. 5.1
AES-256			
Fault Model	No. of Fault	Exhaustive search size	Reference
1 byte random fault	2	$2^{16}$	Sec. 5.2
1 byte random fault	3	1	Sec. 5.3

### 5.1 DFA Attack Against AES-192 based on 2 Faults

The proposed DFA attack against AES-192 is based on two fault injections, one is injected at the beginning of the 10-th round (indexed as  $a$ ) and the other is at the beginning of the 9-th round (indexed as  $b$ ). The propagation of the injected faults is shown in Fig. 2. We denote the state as the input or the output of a certain operation. For example,  $SB_{a12}^{in}$  and  $SB_{a12}^{out}$  are the input state and the out state for the SubBytes in the 12-th round with the fault injected at 10-th round. The requirement is that only one byte of the intermediate value is disturbed, and attackers know which byte it is but does not know the value of the injected fault.



**Fig. 2.** Fault propagation in the DFA attack in AES-192

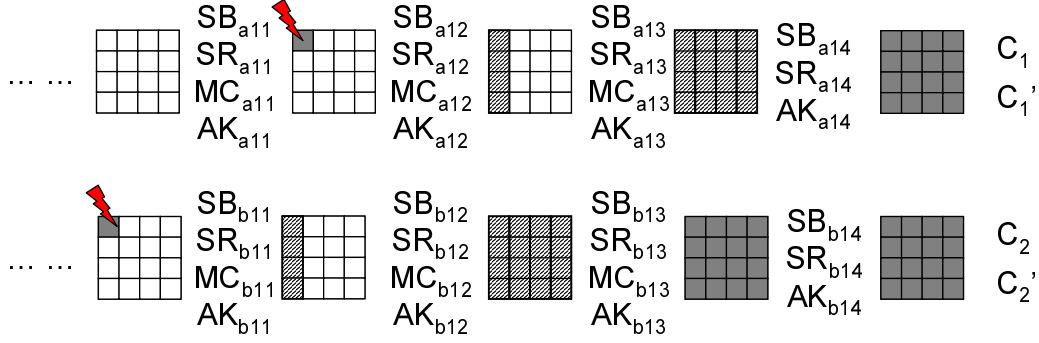
This attack procedure can identify the key without any exhaustive search involving plaintext. Here is the attack procedure:

- Step 1** Restrict the key space of  $K_{12}$  to  $2^{32}$  by checking the difference at each column of the state  $SB_{a12}^{in}$ .
- Step 2** Take a possible candidate of  $K_{12}$ , calculate the corresponding left half of  $K_{11}$ . Use the known part of  $K_{11}$  to calculate partial of the difference at state  $SB_{a11}^{in}$  and  $SB_{b11}^{in}$ . Check whether or not the calculated byte-wise differences satisfy the relationship decided by MixColumns calculation.
- Step 3** Repeat Step 2 for all possible  $K_{12}$ , then  $K_{12}$  is expected to be identified.
- Step 4** Use the identified  $K_{12}$  and the identified left half of  $K_{11}$  to identify the differential at  $SB_{b11}^{in}$  and to calculate the differential at state  $SB_{b11}^{out}$ . Use the S-box differential distribution table to obtain the possible real values at  $SB_{b11}^{out}$ . Use these values to calculate the possible candidates of the right half of  $K_{11}$ . The size of the key space is expected to be  $2^8$ .
- Step 5** For each key guess of  $K_{11}$ , calculate and verify the differentials at state  $SB_{a11}^{in}$  and  $SB_{b10}^{in}$ , then  $K_{11}$  can be identified as well.

The technique used in Step 4 can also be applied to Step 1 to speed up the calculation by omitting the exhaustive search with  $2^{32}$  keys. This attack procedure has been verified by simulating it using the C language. Based on a PC with a Pentium (R) 3.2 GHz CPU and 3.5G RAM, the 192-bit secret key can be identified in about 1 minute without the information of the plaintext.

## 5.2 DFA Attack Against AES-256 based on 2 Faults

The proposed DFA attack against AES-256 is also based on two fault injections, one at the beginning of the 12-th round and the other one at the beginning of the 11-



**Fig. 3.** Fault propagation in the DFA attack in AES-256

th round. The propagation of the injected faults is shown in Fig. 3. The requirement is that only one byte of the intermediate value is disturbed, while attackers know which byte it is but does not know the value of the injected fault.

Here is the attack procedure.

- Step 1** Restrict the key space of  $K_{14}$  to  $2^{32}$  by checking the difference at each column of the state  $SB_{a14}^{in}$ .
- Step 2** Take a possible candidate of  $K_{14}$ , calculate the rightmost three columns of  $K_{12}$  and the differences at state  $SB_{b13}^{out}$ . For the difference of each column at  $SB_{b13}^{in}$ , there are  $2^8$  candidates. Use the S-box differential table to obtain all the possible real values for each column at  $SB_{b13}^{in}$ . After that, the real values for each column at  $SB_{b13}^{in}$  is expected to have  $2^8$  candidates.
- Step 3** Use the the rightmost three columns of  $K_{12}$  and the real values at state  $SB_{b13}^{in}$  to calculate three active bytes at  $SB_{b12}^{in}$ . Check whether or not the calculated active bytes satisfy the relationship of the MixColumns calculation. After this calculation, the space of real values of the rightmost three columns at  $SB_{b13}^{in}$  can be reduced to  $2^8$ .
- Step 4** One can use the possible real values at  $SB_{b13}^{in}$  and the corresponding ciphertext  $C_2$  to calculate the candidates of  $K^{13}$ , the space of  $K^{13}$  is expected to be  $2^{8+8} = 2^{16}$ . Then one can test each candidates by calculating and checking the difference at  $SB_{a13}^{in}$  and  $SB_{b12}^{in}$ . The probability of satisfying all the remaining requirements is about  $2^{-24-8} = 2^{-32}$ , where  $2^{-24}$  for the difference at  $SB_{a13}^{in}$  and  $2^{-8}$  for the difference at  $SB_{b12}^{in}$ .
- Step 5** Repeat Step 1 to Step 4 for all the candidates of  $K^{14}$ , finally the key space of the 256-bit key can be restricted to  $2^{16}$ . Then one can apply exhaustive search to identify the secret key.

With the most used fault model, this procedure can identify the 256-bit secret key with 2 fault injections within a practical calculation complexity ( $2^{48}$ ). The attack complexity can be reduced if two faulty ciphertexts are obtained from the same plaintext, since attackers can restrict the space of  $K_{13}$  to  $2^8$  instead of  $2^{16}$  before Step 4, then the attack complexity can be reduced to  $2^{40}$ .

Based on the usage of the differential table of S-box (well used in the rebound attacks [16]), this attack procedure is able to check the differences at  $SB_{b12}^{in}$  without any knowledge of  $K_{13}$ . This improvement is the key point making the DFA attacks on AES-256 with two fault injections practical.

### 5.3 DFA Attack Against AES-256 based on 3 Faults

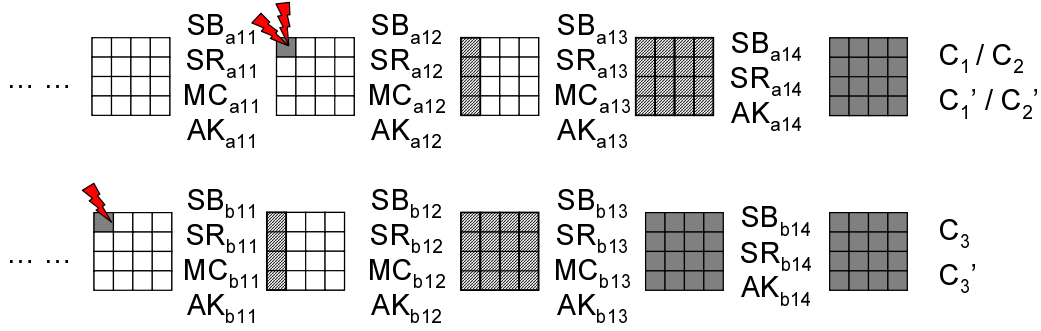


Fig. 4. Fault propagation in the DFA attack in AES-256

This section shows how to identify the 256-bit secret key with 3 fault injections without any exhaustive search based on plaintext. This attack requires two 1-byte fault injections at the beginning of the 12-th round and another 1-byte fault injection at the beginning of the 11-th round. The propagation of the injected faults is shown in Fig. 4. Only one byte of the intermediate value is disturbed, while attackers know which byte it is but does not know the value of the injected fault.

Here is the attack procedure.

**Step 1** Restrict the key space of  $K_{14}$  to  $2^{32}$  by checking the difference at each column of the state  $SB_{a14}^{in}$  based on  $(C_1, C'_1)$ . Then identify  $K_{14}$  by checking the difference at each column at  $SB_{a14}^{in}$  based on  $(C_2, C'_2)$ .

**Step 2** Restrict the key space of  $K_{13}$  to  $2^{32}$  by checking the difference at each column of the state  $SB_{b13}^{in}$  based on the identified  $K_{14}$  and  $(C_3, C'_3)$ .

**Step 3** For the identified  $K_{14}$  and each candidate of  $K_{13}$ , check the difference at state  $SB_{b12}^{in}$  based on  $(C_3, C'_3)$  and the difference at  $SB_{a13}^{in}$  based on both  $(C_1, C'_1)$  and  $(C_2, C'_2)$ . After step 3, the correct key is expected to be identified.

Compared to the introduced DFA attack on AES-256 in Sec. 5.2, this attack takes one more fault injection, but costs much less computation in the key retrieval process and be able to identify the key without knowing the value of plaintext.

## 6 Future Research of DFA against AES

Notice that in the DFA attack flow we proposed, DFA attacks first use the information of  $\Delta I$  at the target state to restrict the key space. Then for the restricted key space, the  $\Delta I$  at injection state is calculated to restrict the key space again. Finally, an exhaustive search based on  $(P, C)$  is applied to identify the key. Divide and conquer makes the information of  $\Delta I$  at the target state can be used in a practical time. The exhaustive search based on  $(P, C)$  is quiet difficult to be further improved. A possible future work for DFA attacks against AES is to find a method to speed up the DFA search up to the injection state.

Assume that attackers get a pair of  $(P, C)$  and get only one faulty ciphertext  $C'$  with a fault that is injected trying to follow the Piret's fault model. When the injected fault actually belongs to Piret's fault model, the key can be fully retrieved in a practical time. Otherwise, we can consider that multiple faulty bytes rather than a single faulty byte is disturbed. If the multiple faulty bytes locate at a diagonal of  $\Delta I$  at the injection state, the key can also be fully retrieved in a practical time [15] as well.

However, the injected fault could be two faulty bytes locate at two diagonals at the beginning of the 8<sup>th</sup> round. In this case, after exploiting the first type of information about  $\Delta I$  at the target state, the key space can be restricted to  $2^{74.3}$ . After that, exploiting the second type of information about  $\Delta I$  at the target state can restrict the key space to  $2^{66.54}$ . Then, theoretically, the third type of information about  $\Delta I$  at the injection state can be used to restricted the key space to  $2^{22.57}$  that is small enough for an exhaustive search based on  $(P, C)$ . Unfortunately, the technical used in Sec. 5 seems difficult to be applied here. Finding a method to exploit the third type of information in a practical time could be an interesting future work <sup>1</sup>.

---

<sup>1</sup> The size of key space is calculated based on the proposed prediction method.



## 7 Conclusions

This paper analyzes differential fault attacks against AES from an information theoretic perspective. The assumptions of fault injection are reviewed as the information of two intermediate values. DFA attacks against AES are considered as the cryptanalysis that obtains the information of key based on the two ciphertexts and the information of the fault injection. Several previous DFA works were reviewed from the information theoretic perspective. Based on our analysis, every fault model has a limitation of the attack efficiency and we proposed a method to predict the attack efficiency for all similar DFA attacks. We gave a general DFA attack flow which requires the least faulty ciphertexts with a reasonable computational cost. We also managed to propose practical DFA attacks against AES-192 and AES-256 reaching the theoretically minimal fault injection times.

## 8 Acknowledgement

The authors would like to thank Chong Hee Kim for kindly showing us his latest research results and his valuable comments to the contents in Sec. 5.

## References

1. W. Li, D. Gu, Y. Wang, J. Li, and Z. Liu. "An extension of differential fault analysis on AES." in *International Conference on Network and System Security*, pages 443–446. IEEE Computer Society 2009.
2. J. Takahashi and T. Fukunaga. "Differential fault analysis on AES with 192 and 256-bit key." Cryptology ePrint Archive, Report 2010/023, 2010. <http://eprint.iacr.org/>
3. A. Barenghi, G. Bertoni, L. Breveglieri, M. Pelliccioli, and G. Pelosi. "Low voltage fault attacks to AES and RSA on general purpose processors." Cryptology ePrint Archive, Report 2010/130, 2010. <http://eprint.iacr.org/>
4. C. H. Kim. "Differential Fault Analysis against AES-192 and AES-256 with Minimal Faults." to appear in *FDTC*, Springer 2010.
5. E. Biham and A. Shamir. "Differential Fault Analysis of Secret Key Cryptosystems." in *CRYPTO*, pages 513–525. Springer 1997.
6. D. Saha, D. Mukhopadhyay and D. RoyChowdhury, "A Diagonal Fault Attack on the Advanced Encryption Standard." Cryptology ePrint Archive, Report 2009/581, 2009. <http://eprint.iacr.org/>
7. D. Boneh, R. A. DeMillo, and R. J. Lipton. "On the Importance of Checking Cryptographic Protocols for Faults. (Extended Abstract)" in *EUROCRYPT*, pages 37–51. Springer, 1997.
8. J. Blömer and J.-P. Seifert. "Fault based Cryptanalysis of the Advanced Encryption Standard." Cryptology ePrint Archive, Report 2002/075, 2002. <http://eprint.iacr.org/>
9. C. Giraud. "DFA on AES." Cryptology ePrint Archive, Report 2003/008, 2003. <http://eprint.iacr.org/>
10. S. Gomisawa, M. Izumi, Y. Li, J. Takahashi, T. Fukunaga, Y. Sasaki, K. Sakiyama, and K. Ohta. "Applicability Extension and Efficiency Improvement of Fault Analysis Attack on AES Implementations." In *SCIS*, 2010.
11. A. Moradi, M. T. Manzuri Shalmani, and M. Salmasizadeh. "A Generalized Method of Differential Fault Attack against AES cryptosystem." In *CHES*, pages 91–100, 2006.

12. D. Mukhopadhyay. “An Improved Fault based Attack of the Advanced Encryption Standard.” In *AFRICACRYPT*, pages 421–434, 2009.
13. G. Letourneux P. Dusart and O. Vivolo. “Differential Fault Analysis on A.E.S.” Cryptology ePrint Archive, Report 2003/010, 2003. <http://eprint.iacr.org/>
14. G. Piret and J.-J. Quisquater. “A Differential Fault Attack Technique against SPN structures, with Application to the AES and  $K_{HAZAD}$ .” In *CHES*, pages 77–88, 2003.
15. M. Tunstall and D. Mukhopadhyay. “Differential Fault Analysis of the Advanced Encryption Standard using a Single Fault.” Cryptology ePrint Archive, Report 2009/575, 2009. <http://eprint.iacr.org/>
16. F. Mendel and C. Rechberger and M. Schl affer and S. S. Thomsen “The Rebound Attack: Cryptanalysis of Reduced Whirlpool and Gr ostl.” In *FSE*, pages 260–276, 2009.