

Bricklayer Attack: A Side-Channel Analysis on the ChaCha Quarter Round

Alexandre Adomnicai^{1,3}, Jacques J. A. Fournier²,
and Laurent Masson¹

¹ Trusted Objects, Aix-en-Provence, France
{a.adomnicai,l.masson}@trusted-objects.com

² CEA-Leti, Grenoble, France

jacques.fournier@cea.fr

³ EMSE, Gardanne, France

Abstract. ChaCha is a family of stream ciphers that are very efficient on constrained platforms. In this paper, we present electromagnetic side-channel analyses for two different software implementations of ChaCha20 on a 32-bit architecture: one compiled and another one directly written in assembly. On the device under test, practical experiments show that they have different levels of resistance to side-channel attacks. For the most leakage-resilient implementation, an analysis of the *whole* quarter round is required. To overcome this complication, we introduce an optimized attack based on a divide-and-conquer strategy named *bricklayer attack*.

Keywords: ChaCha · Implementation · Side-channel attacks

1 Introduction

ChaCha [7] is a family of stream ciphers introduced by Daniel J. Bernstein in 2008. It is a variant of the Salsa20 family [8], which is part of the eSTREAM portfolio [4], providing better diffusion for similar performances. ChaCha is an ARX-based cipher, which means that it only uses modular additions, rotations and bitwise XORs. It has been widely adopted for encryption, as well as for random number generation in many operating systems (*e.g.* Linux, OpenBSD) and protocols (*e.g.* SSH, TLS). Moreover, the upcoming version 1.3 of the Transport Layer Security (TLS) protocol [35] will allow Authenticated Encryption with Associated Data (AEAD) cipher suites only, leaving AES-CCM [31], AES-GCM [37] and ChaCha20-Poly1305 [25] as the only three options. This update should significantly increase the use of ChaCha in the near future. On top of that, the Internet of Things (IoT) should be in favour of the ChaCha deployment (*e.g.* Apple HomeKit for IoT devices [2]), since its instances are cheaper than AES on microcontrollers that do not have any dedicated cryptographic hardware. For instance, on Android phones, HTTPS connections from Chrome browsers to Google now use ChaCha20-Poly1305 [12].

As a result of its standardization, ChaCha is under close scrutiny with regards to cryptanalysis, especially regarding differential attacks [3, 14, 28, 38, 40]. Recently, studies have been carried out to evaluate its physical security, especially regarding fault attacks [24, 32]. However, only one side-channel analysis has been proposed so far [21]. We believe that further work must be undertaken in this field since ChaCha is particularly well suited for embedded devices.

Our Contribution. In this paper, we focus on the side-channel analysis of ChaCha by taking two different implementations into consideration.

First, we investigate the OpenSSL C source code compiled on a 32-bit ARM microcontroller. It results in a straightforward attack path, which consists in targeting each 32-bit key word independently.

The second target is an assembly implementation which saves some memory accesses. We highlight that, on the device under test (DUT), this slight modification protects from the only side-channel attack published to date. Nevertheless, our implementation remains vulnerable even though attack paths are more complex. We tackle this problem by introducing the *bricklayer attack*, which is based on a divide-and-conquer approach, and emphasize that attacking from the keystream rather than from the input is way more efficient.

Outline. First, we present the ChaCha family of stream ciphers before providing an outline of side-channel attacks. Then, we describe our approaches on performing electromagnetic analyses depending on software implementations of ChaCha. Subsequently, we present our practical results and discuss the feasibility of conducting these attacks in real-world scenarios. Finally, we analyze the overhead introduced by the masking countermeasure in the specific case of ChaCha20.

2 The ChaCha Family of Stream Cipher

As its predecessor, and unlike traditional stream ciphers, ChaCha does not have an initialization phase since it works like a block cipher used in counter (CTR) mode [18]. Its core is an ARX-based function which maps a 512-bit input block to a 512-bit output key stream. Input blocks are built by arranging data in a 4×4 matrix where each element is a 32-bit word. The encryption key fills half of the matrix as it is 256-bit long, while the two remaining quarters are respectively occupied by the inputs and the constant ‘expand 32-byte k’. This constant aims at reducing the amount of data an attacker can control while the inputs refer to a nonce which is built from the block counter and the initial vector (IV) (Fig. 1).

The core function is defined by iterating several rounds on the input block, where each round consists of four parallel quarter round (QR) operations. A QR updates 4 words (*i.e.* a block quarter) as defined in Algorithm 1 where \boxplus means addition modulo 2^{32} , \oplus means XOR and \lll means left bitwise rotation.

Depending on the round number (enumerated from 0), each QR operates either on a column, or on a diagonal. ChaChaR refers to a specific instance

‘expa’	‘nd 3’	‘2-by’	‘te k’
k_0	k_1	k_2	k_3
k_4	k_5	k_6	k_7
nonce ₀	nonce ₁	nonce ₂	nonce ₃

Fig. 1. ChaCha’s input block initialization

Algorithm 1. ChaCha quarterround(a , b , c , d)

$a \boxplus= b;$	$d \oplus= a;$	$d \lll= 16;$
$c \boxplus= d;$	$b \oplus= c;$	$b \lll= 12;$
$a \boxplus= b;$	$d \oplus= a;$	$d \lll= 8;$
$c \boxplus= d;$	$b \oplus= c;$	$b \lll= 7;$

where R rounds are used. Several variants are defined with 8, 12 or 20 rounds, defining different trade-offs between security and performance. Recently, it has been shown under certain assumptions that ChaCha12 is sufficiently secure to ensure a 256-bit security level [14]. Nevertheless, ChaCha20 remains the most widespread instance for security margins. In many implementations, ChaCha R uses $\frac{R}{2}$ iterations of double rounds instead of R rounds, which consists in a column round and a diagonal one.



On top of iterating several rounds on the input block, an additional step is required. The reason is that while QRs scramble blocks beyond recognition, they are invertible. Therefore, applying the reverse of each operation in the reverse order leads to the original block and thus, the encryption key. ChaCha prevents this by adding the original block to the scrambled one, word by word, in order to generate the pseudo-random block. The whole encryption process is detailed in Algorithm 2.

3 Background on Side-Channel Attacks

3.1 Correlation Electromagnetic Analysis

Cryptographic primitives are usually built to resist to mathematical cryptanalysis or exhaustive key search. However, they are designed to be finally executed

Algorithm 2. ChaChaR encryption

Require:

n -bit plaintext P
 encryption key k
 counter ctr
 IV iv

Ensure: n -bit ciphertext C **for** i from 0 to $\lfloor n/512 \rfloor$ **do** $B \leftarrow \text{init}(k, ctr, iv)$ ▷ input block initialization $B' \leftarrow B$ ▷ working variable**for** j from 0 to $\frac{R}{2} - 1$ **do**quarterround($B'_0, B'_4, B'_8, B'_{12}$) ▷ column roundsquarterround($B'_1, B'_5, B'_9, B'_{13}$)quarterround($B'_2, B'_6, B'_{10}, B'_{14}$)quarterround($B'_3, B'_7, B'_{11}, B'_{15}$)quarterround($B'_0, B'_5, B'_{10}, B'_{15}$) ▷ diagonal roundsquarterround($B'_1, B'_6, B'_{11}, B'_{12}$)quarterround($B'_2, B'_7, B'_8, B'_{13}$)quarterround($B'_3, B'_4, B'_9, B'_{14}$)**end for** $B \leftarrow B \boxplus B'$ ▷ final block addition $C_i \leftarrow P_i \oplus B$ $ctr \leftarrow ctr + 1$ **end for**

on a given processor with its own physical characteristics. Electronic circuits are inherently leaky as they produce emissions that make it possible for an attacker to deduce how the circuit works and what data is being processed. Because these emissions are nothing more than side effects, their use to recover cryptographic keys has been termed ‘side-channel attacks’. Since the publication of Differential Power Analysis (DPA) [23], it is common knowledge that the analysis of the power consumed by the execution of a cryptographic primitive might reveal information about the secret involved.

A few years later, Correlation Power Analysis (CPA) has been widely adopted over DPA as it requires fewer traces and has been shown to be more efficient [11]. The principle is to target a sensitive intermediate state of the algorithm and try to predict its value from the known input and different key guesses. Then, to uncover the link between these predictions and the leakage measurements, the Pearson correlation coefficient between these two variables is computed using an appropriate leakage model. The Hamming weight (HW) and the Hamming distance (HD) model are the most commonly used models to simulate the leakage of a cryptographic device. For each key hypothesis, it results in a value between -1 (total negative correlation) and 1 (total positive correlation) for every point in time, indicating how much the prediction correlates with the recorded values over several measurements. The formula of this coefficient is

$$\text{Corr}(X, Y) = \frac{E(X \cdot Y) - E(X) \cdot E(Y)}{\sqrt{E((X - E(X))^2) \cdot E((Y - E(Y))^2)}} \quad (1)$$

where $E(X)$ is the expected value of the random variable X . Finally, the hypothesis which matches with the real key should return a significantly higher coefficient than the other hypotheses. This attack remains valid when analyzing electromagnetic emanations [19,34] instead of power consumption, since they are mainly due to the displacement of current through the rails of the metal layers. In this case, we talk about Correlation ElectroMagnetic Analysis (CEMA).

3.2 Selection Function

The intermediate state y on which the side-channel attack focuses is defined by a *selection function* $\varphi(x, k) = y$, which is part of the encryption algorithm. It depends on x , a known part of the input and on k , an unknown part of the secret key. Usually, selection functions are chosen to be easy to compute, typically at the beginning of the encryption or decryption process. Furthermore, a valuable property for selection functions is high non-linearity as it ensures a good distinguishability between the correct and incorrect key guesses. Indeed, correlation between the leakage and the prediction will be close to zero if the key guess is incorrect due to their non-linear relationship.

In case of ARX structures, the non-linearity only relies on modular additions, while diffusion is provided by rotations (diffusion within single words) and XORs (diffusion between words). Although the carry propagation in the modular addition results in some non-linearity, it is not as good a candidate as S-boxes. It can be explained by the fact that most significant bits in the output of a modular addition are more subject to non-linearity than least significant ones. However, side channel attacks remain possible as shown in numerous publications [10,26,41].

4 Side-Channel Overview of ChaCha

4.1 ChaCha Case Study

To set up such a side-channel attack, one has to determine an attack path (*i.e.* to choose a selection function) either starting from the plaintext, or from the ciphertext. Physical attacks against stream ciphers can be challenging because the key stream is computed independently from the plaintext/ciphertext, which interferes in the relationship between known values and the secret key. However, from a side-channel point of view, ChaCha differs significantly from other stream ciphers' designs such as linear-feedback shift registers where the key is only directly involved during registers' initialization. Indeed, as ChaCha operates like a block cipher in CTR mode, the key is directly manipulated everytime a 512-bit block needs to be encrypted. More precisely, each key word directly interacts

with other data during the first round (after which they have been updated) and again during the final block addition.

An attack that takes advantage of the first round has already been published in [21]. The attack on the i^{th} column round ($0 < i < 4$) relies on the selection function defined by

$$\varphi_0(\text{nonce}_i, \tilde{k}_i \parallel k_{i+4}) = \left((\text{nonce}_i \oplus \tilde{k}_i) \lll 16 \right) \boxplus k_{i+4} \quad (2)$$

where $\tilde{k}_i = k_i \boxplus \text{constant}_i$. However, this selection function forces the attacker to target two key words at once, which results in a key search space $|\mathcal{K}| = 2^{64}$. Since the bit-size of the targeted subkey determines the memory complexity of the side-channel attack, one can understand why this would be undoable in practice. To get around this problem, the authors exploit the QRs' intermediate states in order to operate step by step. They propose to first recover k_i by targeting $\text{nonce}_i \oplus \tilde{k}_i$ and then take advantage of its knowledge to find k_{i+4} . Therefore, recovering k_i and k_{i+4} requires the knowledge of nonce_i . However, the paper also describes an attack path that allows to recover the entire key with the knowledge of only two words. This latter exploits several intermediate states in the first two rounds.

Regarding the final block addition, an attacker could choose $\varphi(x, k) = x \boxminus k$ where x refers to a keystream word and \boxminus refers to modular subtraction. Compared to the previous attack path, it has the advantage of recovering all key words using the modular subtraction as selection function. Moreover, all keystream words are pseudorandom values, which is not necessarily the case for nonces. However, this selection function requires the knowledge of the keystream (*i.e.* both plaintext and ciphertext).

Throughout this paper, we will make the assumption that an attacker has access to all this information. In Sect. 6 we discuss the attacks' feasibility in practice and thus, whether our assumptions are reasonable.

4.2 Implementation Aspects

When targeting software implementations on load/store architectures, data transfers due to memory accesses (*i.e.* loads and stores between memory and registers) are known to leak the most information compared to arithmetic and logic operations [13, 30], which only occur between registers and are usually unexploitable in practice [9]. Our practical experiments on the DUT presented in Sect. 6 verified this hypothesis. Therefore, the intermediate values that are manipulated by these sensitive operations should be easiest to target, introducing a direct link between selection functions and implementation aspects.

Throughout this paper we will study selection functions in relation to memory accesses, assuming they are the main source of exploitable leakage.

4.3 OpenSSL Implementation

First, we decided to attack a C implementation of ChaCha20 in order to see how compilers can deal with ARX structures and memory accesses. To do so,

we compiled the ChaCha20 C implementation from OpenSSL (version 1.0.1f) for an ARM Cortex-M3 microcontroller using the GNU ARM C compiler 5.06 (update 2). Regardless of the optimization level chosen (from `-O0` to `-O3`), within a QR, each addition and each rotation is followed by a STR instruction. Hence, these memory accesses allowed us to carry out the attacks described above. Practical results are briefly presented in Sect. 6 for comparative purposes.

4.4 Side-Channel Analysis of the Salsa20 Quarter Round

In the next section, we show how memory accesses can be easily managed to remove the leakage of intermediate states within a QR. This implies to target the QR output without taking its intermediate values into consideration, making the attacks presented in [21] irrelevant in this case. Although such an analysis has already been performed on Salsa20 [29], it does not apply to ChaCha.

Algorithm 3. Salsa20 quarterround(a, b, c, d)

$$\begin{array}{ll} b \oplus = (a \boxplus d) \lll 7; & c \oplus = (b \boxplus a) \lll 9; \\ d \oplus = (c \boxplus b) \lll 13; & a \oplus = (d \boxplus c) \lll 18; \end{array}$$

In the case of Salsa20, as described in Algorithm 3, the update of the second input only depends on itself and two others (the first and the last). This allows to recover the key words involved in this computation as first/last input words, with two other ‘non-key’ operands (*i.e.* constant and nonce). The attack consists in performing a CPA on a 32-bit value using a divide-and-conquer (D&C) approach, which consists in separating the attack into $\lceil \frac{32}{n} \rceil$ computations on n -bit windows in parallel. The other key words that do not match these requirements were retrieved by using the knowledge of those which have been previously recovered. This allowed to keep a search space of 2^{32} instead of 2^{64} . On top of providing better diffusion, the ChaCha QR gives *each* input word a chance to affect the other three twice. This adjustment makes the attack irrelevant against ChaCha since the key search space cannot be less than 2^{64} in any case.

5 Side-Channel Analysis of the Quarter Round

Throughout this section, for greater clarity, we assume that all operators are left-associative so that

$$a \boxplus b \oplus c \lll d \iff (((a \boxplus b) \oplus c) \lll d).$$

5.1 Optimizing Memory Accesses

A solution to overcome attacks on intermediate states within QRs is a straightforward assembly implementation, which is a good way to reduce memory access instructions for load/store architectures. As explained in [9], for some instances of ARX lightweight block ciphers like Simon and Speck [5], it is possible to keep the whole state in registers during the entire encryption process. Thereby, they can be implemented in assembly without having to execute a single STR instruction during the whole encryption process, drastically reducing the amount of leakage.

Unfortunately, in the case of ChaCha, the state consists of 16 32-bit words. Therefore, it would require a 32-bit CPU with at least 16 general-purpose registers (excluding the stack pointer, the program counter and other specific cases such as hardwired registers) to avoid memory accesses. As our chip only has 13 general-purpose registers, we implemented ChaCha so that word values are loaded into registers at the beginning of each QR and are then stored in RAM at the end. Furthermore, during the last round, related key words are also loaded into registers at the beginning of QRs, resulting in

$$\begin{aligned} &\text{quarterround}'(x_0, x_5, x_{10}, x_{15}, k_1, k_6) \\ &\text{quarterround}'(x_1, x_6, x_{11}, x_{12}, k_2, k_7) \\ &\text{quarterround}'(x_2, x_7, x_8, x_{13}, k_3, k_4) \\ &\text{quarterround}'(x_3, x_4, x_9, x_{14}, k_0, k_5) \end{aligned}$$

where $\text{quarterround}'(a, b, c, d, x, y) = \text{quarterround}(a, b, c, d) \boxplus (0, x, y, 0)$. This method protects against leakages that would allow an attack from the keystream using the modular subtraction as selection function. Thus, these elementary implementation tricks imply to analyze the side-channel resilience of the *whole* QR.

5.2 Focusing on the Quarter Round

As every word influences the three others, and is updated twice, the simplest selection function would be defined by focusing, during the first column rounds, on the word which is completely updated at first, resulting in having

$$\varphi_1(\text{nonce}_i, k_i \parallel k_{i+4}) = \text{nonce}_i \oplus \tilde{k}_i \lll 16 \boxplus k_{i+4} \oplus k_i \lll 12 \boxplus \tilde{k}_i. \quad (3)$$

However, as previously mentioned, this implies a side-channel attack on 64 bits, which is not feasible in practice. Therefore, we investigated the relevance of the D&C approach in this specific case. Figure 2 sketches how key words are involved in computations. It results that targeting n bits of $y = \varphi_1(\text{nonce}_i, k_i \parallel k_{i+4})$ does not lead to a complexity equal to 2^{2n} since rotations make different n -bit windows interact with each other. As there is a rotation of 16 bits followed by another one of 12, some bits of \tilde{k}_i may overlap. Hence, the key search space depends on

the windows' size.

$$|\mathcal{K}| = \begin{cases} 2^{4n}, & \text{if } n \leq 4 \\ 2^{3n+4}, & \text{if } 4 \leq n \leq 12 \\ 2^{2n+16}, & \text{if } 13 \leq n \leq 16 \\ 2^{n+32}, & \text{otherwise} \end{cases} \quad (4)$$

Furthermore, rotations are discarded from the selection function, resulting in

$$\varphi_{2,n}(\text{nonce}_i, \tilde{k}_i^A \parallel k_i^B \parallel k_{i+4}^B \parallel \tilde{k}_i^C) = \text{nonce}_i^A \oplus \tilde{k}_i^A \boxplus_n k_{i+4}^B \oplus k_i^B \boxplus_n \tilde{k}_i^C \quad (5)$$

where superscripts refer to intervals that define n -bit windows.

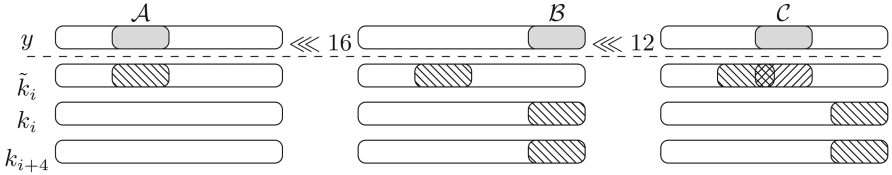


Fig. 2. D&C approach on the ChaCha QR, $n = 8$

In order to evaluate this method, we performed software simulations using the HW model (without any additional noise) and random nonces. As expected, the right key matches with the highest correlation coefficient. Nevertheless, some other hypotheses also lead to the maximum coefficient as shown in Fig. 3, resulting in collisions.

Definition 1 (Collision). Let $\varphi(n, k)$ be a selection function and κ be the right key hypothesis. A collision is an hypothesis κ' such that $\varphi(n, \kappa) = \varphi(n, \kappa')$ for all n .

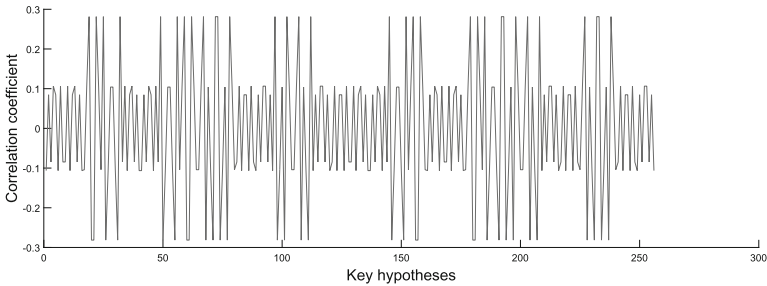


Fig. 3. Attack simulation on $\varphi_{2,2}$

Proposition 1. *An attack on $\varphi_{2,n}$ returns up to $n \cdot 2^{n+2}$ collisions.*

Another point that has not been discussed so far is the drawback caused by carry propagations. Except when focusing on the least significant bits (LSBs), one has no way of knowing if subkeys involved in additions are affected by a carry. Thus, the positions of targeted windows are very important. Plus, we made the choice to dissociate \tilde{k}_i from k_i in order to prevent from erroneous predictions of $k_i^A \boxplus_n \text{constant}_i^A$ and $k_i^C \boxplus_n \text{constant}_i^C$. For instance, in Fig. 2, \tilde{k}_i^C is the only hypothesis which could be erroneous due to a carry propagation on its addend. As a result, an attacker should mount one attack taking this carry into consideration, and another one without. This would mean that the total number of collisions would be doubled. Although this selection function may provide some information, we chose to investigate a more efficient attack path.

5.3 Benefits of the Reverse Function

The ChaCha QR is trivially invertible and the inverse quarter round (IQR) is defined in Algorithm 4.

Algorithm 4. ChaCha inv_quarterround(a, b, c, d)

$b \ggg 7;$	$b \oplus= c;$	$c \boxminus= d;$
$d \ggg 8;$	$d \oplus= a;$	$a \boxminus= b;$
$b \ggg 12;$	$b \oplus= c;$	$c \boxminus= d;$
$d \ggg 16;$	$d \oplus= a;$	$a \boxminus= b;$

What matters here is that each input word does not have a chance to influence the other three, since the first word does not impact the update of the second one. Hence, the overall selection can be defined as below

$$\varphi_3(b \parallel c \parallel \tilde{d}_i, k_b \parallel k_c) = (b \boxminus k_b \ggg 7) \oplus (c \boxminus k_c \ggg 12) \oplus (c \boxminus k_c \boxminus \tilde{d}_i) \quad (6)$$

where $\tilde{d}_i = d_i \boxminus \text{nonce}_i$. Regarding the D&C approach where rotations are discarded, it results in the following selection function.

$$\varphi_{4,n}(b \parallel c \parallel \tilde{d}_i, k_b^A \parallel k_c^B \parallel k_c^C) = (b^A \boxminus_n k_b^A) \oplus (c^B \boxminus_n k_c^B) \oplus (c^C \boxminus_n k_c^C \boxminus_n \tilde{d}_i^C) \quad (7)$$

As less words are involved, the key search space is reduced and still depends on the windows' size.

$$|\mathcal{K}| = \begin{cases} 2^{3n}, & \text{if } n \leq 12 \\ 2^{2n+12}, & \text{if } 12 \leq n \leq 20 \\ 2^{n+32}, & \text{otherwise} \end{cases} \quad (8)$$

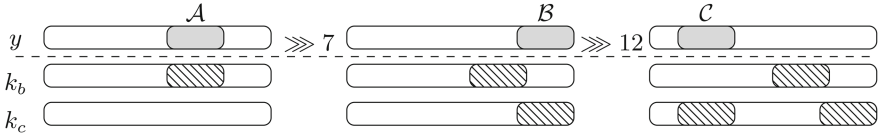


Fig. 4. D&C approach on the ChaCha IQR, $n = 8$

However, since the rotations are less pronounced, key words do not overlap if the windows' size does not exceed 12 bits, as depicted in Fig. 4. Throughout the rest of this section, we only consider the case where $n \leq 12$.

As before, key hypotheses might be affected by carry propagations. However, another advantage of $\varphi_{4,n}$ over $\varphi_{2,n}$ is that one knows the *entire* 32-bit minuend (*i.e.* b or c). Thus, depending on its value, one can calculate the probability of a carry propagation. For instance, when targeting $k_b^{[x,x+n]}$, the probability is

$$p = \mathbb{P}\left(k_b^{[0,x]} > b^{[0,x]}\right) = \frac{2^x - (b^{[0,x]} + 1)}{2^x}. \quad (9)$$

For our simulation with $n = 4$, we took a carry into consideration only if $p > 0.75$. On top of providing a smaller key search space, $\varphi_{4,n}$ is less prone to collisions as shown by our simulation depicted in Fig. 5.

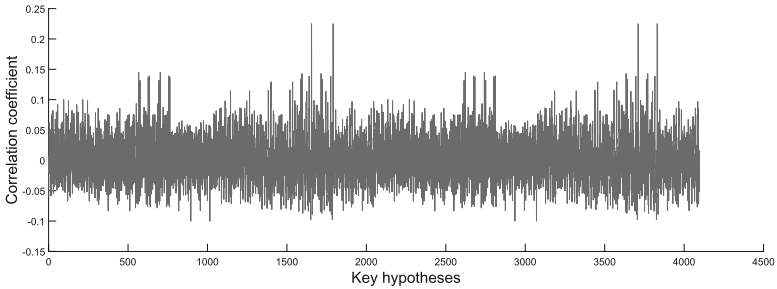


Fig. 5. Attack simulation on $\varphi_{4,4}$

Proposition 2. *An attack on $\varphi_{4,n}$ returns 4 collisions.*

Proof. Flipping the MSB of the minuend/subtrahend also flips the MSB of the modular difference. Therefore, in the case of $\varphi_{4,n}$, flipping the MSB of two n -bit key windows leads to the same output. As a result, the number of collisions is equal to $1 + \binom{3}{2} = 4$. \square

This property allows to halve the key search space (*i.e.* $|\mathcal{K}| = 2^{3n-1}$), since all collisions can be retrieved from just one. In the next section, we suggest a more efficient method than repeating this computation over several windows and then sorting the right key from the collisions.

5.4 Overview of the Brickerlayer Attack

Once collisions have been found using $\varphi_{2,n}$ or $\varphi_{4,n}$, one has to reiterate the same procedure on different windows. Instead of executing several attacks in parallel, we suggest to take advantage of windows that have been previously recovered, in order to target larger ones. For instance, once 4 collisions have been found after an attack on $\varphi_{4,n}$, one can target $\varphi_{4,m}$, where $m > n$, with a complexity $|\mathcal{K}| = 2^{3(m-n)+1}$.

Proceeding in this sequential manner has two advantages. First, taking the carry propagation into consideration is only necessary during the first attack. This property is especially interesting for $\varphi_{2,n}$ since there is no way to estimate carry propagations in this case. Second, each attack cancels collisions from the previous ones, since the positions of the collision bits are changed. For instance, regarding $\varphi_{4,n}$ where collisions only depend on MSBs, the bricklayer approach transforms previous collisions into the predictions' lower bits, allowing the correct collision to stand out. This property is less efficient in the case of $\varphi_{2,n}$ since collisions depends on all bits of the n -bit word. Therefore, the correct collision does not stand out directly but some wrong hypotheses are still discarded.

An example application of the bricklayer attack using $\varphi_{4,n}$ is depicted in Fig. 6. Note that from the fourth step, the attack focuses two key windows instead of three because rotations lead to a position that has already been recovered. Finally, the last step considers the entire 32-bit output word using φ_3 and the known bits/collisions.

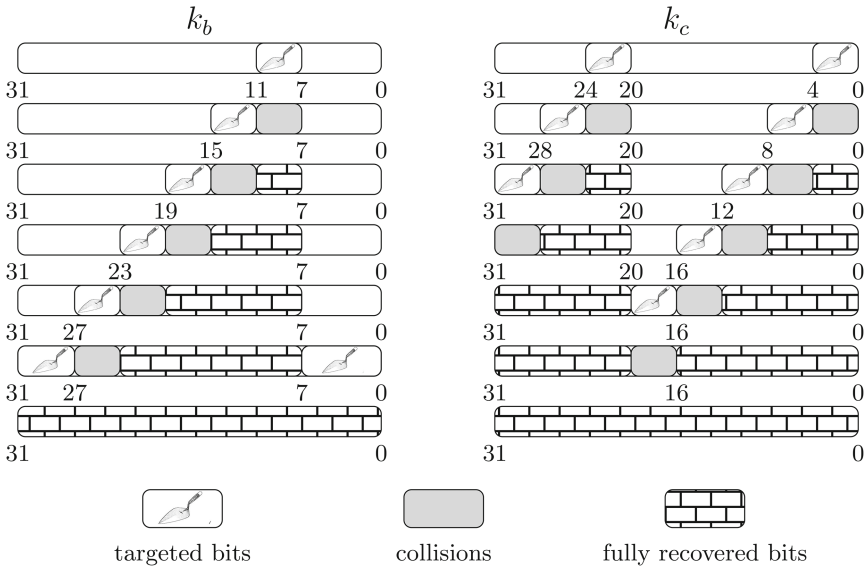


Fig. 6. Bricklayer attack example on IQR

6 Applications in Practice

6.1 Practical Experiments

All practical experiments presented below were done using an ARM 32-bit Cortex-M3 processor clocked at 24 MHz. Note that the DUT does not embed any hardware countermeasure against side-channel attacks. A trigger signal was inserted to indicate the beginning and the end of the penultimate round in order to avoid synchronization complications. EM emanations were measured using a Langer LF-U 5 near-field probe (100 kHz–50 MHz) and a LeCroy WaveSurfer 10 oscilloscope sampled at 10 GS/s. The signal was amplified using a Langer PA 303 BNC preamplifier, providing a gain of 30 dB. We used the same leakage model as for our simulations, since our microcontroller leaks the HW of intermediate values.

First, we tried to perform correlation analyses by focusing on arithmetic operations, without success. Figure 7 emphasizes that attacking the final block addition during executions of `quarterround'` was not successful, whereas for the compiled C version (which stores the intermediate values in RAM), we were able to retrieve the key bits. This reinforced our assumption that, depending on the computing platform, memory accesses can be the only source of exploitable leakage for software implementations.

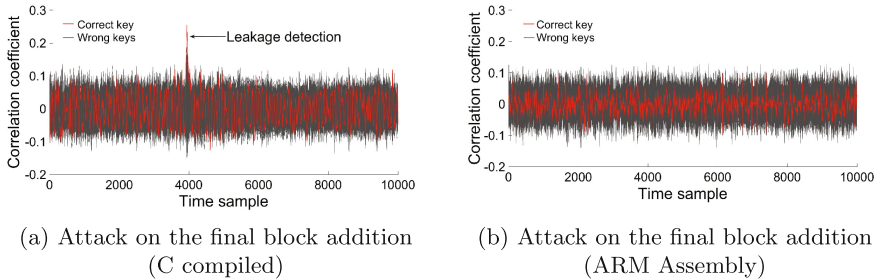


Fig. 7. Impact of memory accesses on electromagnetic leakage

In order to put the bricklayer attack into practice, the following hard-coded input block was used to encrypt 250 kB of data, where the counter (*i.e.* `nonce0`) was incremented for each 512-bit block (Fig. 8).

Figure 9 depicts all the correlation curves corresponding to each step of the bricklayer attack when targeting k_2 and k_7 . We incremented the windows' length by 4 at each step, exactly as illustrated in Fig. 6, resulting in an overall computational complexity of 2^{13} . All CEMAs were computed by halving the key search space. Consequently, some results do not appear clearly on charts and have to be deduced.

The first step, which targets $k_7^{23..20} \parallel k_7^{3..0} \parallel k_2^{10..7}$, returned the collisions $\Gamma = \{\gamma_1, \gamma_2, \gamma_3, \gamma_4\} = \{56, 176, 2096, 2232\}$. For the next stages, each key

61707865	3320646e	79622d32	6b206574
ad0578e5	e962fc0a	42ffc031	75018bee
b7ae69dc	f1490ca8	89ac12fd	be8466d3
00000000	f1d69cbf	8e34191d	7024af3b

Fig. 8. Input block used for practical experiments

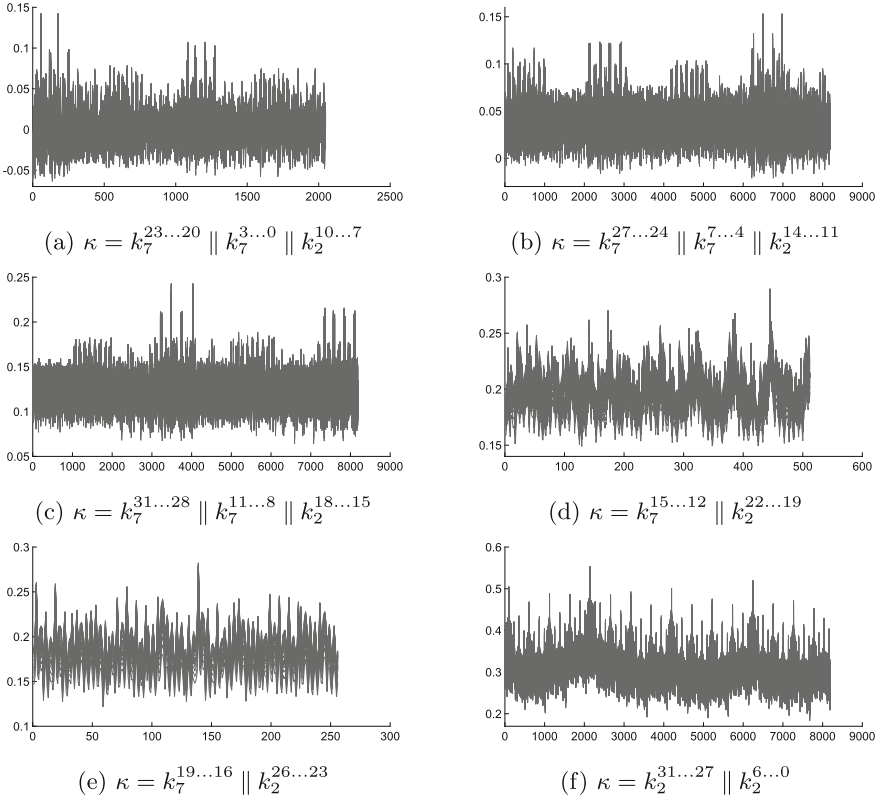


Fig. 9. CEMAs to recover k2 and k7

hypothesis $\kappa \in \mathcal{K}$ was coupled to each collision $\gamma_j \in \Gamma$ and was placed at the index $i = \kappa \cdot |\Gamma| + j$ of the prediction vector. Thus, higher coefficients at indexes i revealed the correct collision of the previous step γ_j by computing $j = i \bmod |\Gamma|$. Finally, the new collisions are equal to $(i - j) / |\Gamma|$. For instance, Fig. 9b indicates that the maximum coefficient appears at indexes $i \in \{6499, 6979\}$. Both indexes are congruent to 3 modulo 4, which means that $\gamma_3 = k_7^{23\dots20} \parallel k_7^{3\dots0} \parallel k_2^{10\dots7}$. As a result, the collisions for $k_7^{27\dots24} \parallel k_7^{7\dots4} \parallel k_2^{14\dots11}$ are defined by $\Gamma = \{1624, 1744, 3664, 3800\}$. The remaining steps followed the same methodology, making it possible to recover k_2 and k_7 entirely. Obviously, this can be applied on other IQRs in parallel to recover the whole encryption key.

A drawback of the D&C method is the number of required measurements, since the leakage of the omitted bits influences the attacked ones. Thus, more traces are needed in order to average out noise. Figure 10 compares, regarding the number of measurements, an attack on the QR using $\varphi_{2,3}$ with the first step of the bricklayer attack presented above, using the same measurement setup.

As a result, to recover the same number of key bits, $\varphi_{4,n}$ requires less traces as it targets larger windows than $\varphi_{2,n}$. However, the number of required traces decreases at each step of the bricklayer attack as the size of targeted windows increases.

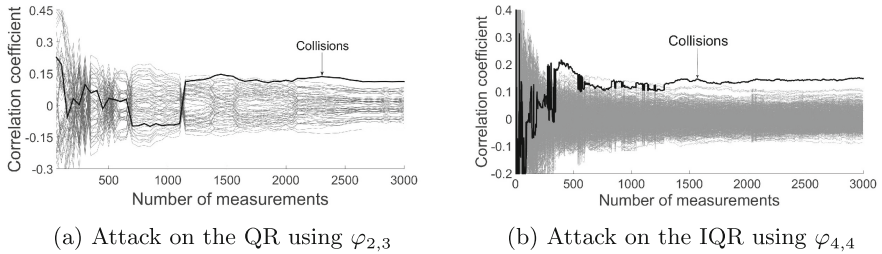


Fig. 10. Correlation coefficients to recover 12 key bits

6.2 Attacks' Feasibility on Existing Protocols

In a typical side-channel analysis, it is assumed that the attacker has access to either the plaintext or the ciphertext, but not necessarily to both. In the case of ChaCha, we can consider the knowledge of the nonce as the knowledge of the plaintext. However, attacks using $\varphi_{4,n}$ require the knowledge of the keystream (*i.e.* plaintexts and ciphertexts), in addition to nonces. This is a strong assumption that could be available in an evaluation laboratory but might be hard to set up in practice, leaving the attacks from nonces more realistic. Therefore, we discuss whether the knowledge of nonces is a fair assumption.

By definition, the single requirement for a cryptographic nonce is to be used only once. Therefore, a simple counter could suit the need. However, in cases where many different keys are used, some protocols (*e.g.* TLS) force a part of the nonce (*e.g.* the IV) to be random in order to thwart multi-key attacks [27]. This leaves the block counter as the only predictable part of the nonce. Therefore, if this latter is defined on n bits, then a correlation analysis cannot recover more than $2 \cdot n$ key bits. As a result, it introduces a protocol-level countermeasure which protects a large part of the key.

Still, existing protocols are not defined in this way. For instance, the Secure Shell (SSH) protocol uses the packet sequence number as a 64-bit IV [1] whereas the remaining 64 bits are used for the block counter, which is reset for each packet. Consequently, observing an entire SSH session makes it possible to predict the entire nonce, giving an attacker the opportunity to recover all key words as soon as enough packets are transmitted.

Furthermore, another construction that can be encountered in practice is XChaCha20, which is implemented in the Sodium crypto library [17]. This construction was first proposed for Salsa20 [6] and aims at extending the nonce to 192 bits so that it can be picked at random. The main idea is to encrypt a block with a fixed key k and 128 bits of the random nonce, without executing the final block addition. The first and last 16 bytes of the output result in a 256-bit subkey k' . Finally, the regular ChaCha20 algorithm is executed using the 64 remaining bits of the 192-bit nonce as IV, and k' as encryption key. Note that XChaCha20 is intrinsically resistant against attacks from the keystream, since the final block addition is omitted during the subkey generation. However, the 192-bit nonce must be transmitted in clear and can be entirely known by the attacker.

These real life case studies introduce the need of dedicated countermeasures against side-channel attacks when ChaCha is deployed in such conditions.

7 Towards a Secure Implementation

A common approach to thwart side-channel attacks is the use of *masking*. This countermeasure consists in blinding the processed values x by means of random masks r , so that intermediate variables are impossible to predict. Thus, an attacker has to analyze multiple point distributions, which exponentially increases the attack complexity with the number of shares. In this section, we only discuss first-order masking *i.e.* the case where a single mask is used to randomize the data. Because of their structures, ARX designs need both boolean ($x' = x \oplus r$) and arithmetic ($x' = x \boxplus r$) masking.

To overcome this complication, there are two main approaches. The first one is to switch from one masking scheme to the other whenever necessary. The first conversion algorithms, described by Goubin in [20], have complexity of $\mathcal{O}(1)$ for boolean to arithmetic and $\mathcal{O}(k)$ for arithmetic to boolean, where k refers to the addends' bit size. The latter was then improved by Coron *et al.* to $\mathcal{O}(\log k)$ [15]. The second approach is to directly perform an addition on the masked values, eliminating the need for conversions [22]. However, secure adders usually rely on the recursion formulae involved in arithmetic to boolean conversions. Consequently, they inherit from the same complexity.

The best method, in terms of performance, depends on the algorithm to be protected. For instance, masks conversions are more efficient when several arithmetic operations are processed successively, since only one arithmetic to boolean conversion is ultimately required. Otherwise, secure adders can lead to better performances as shown by a practical comparison between HMAC-SHA-1 and Speck in [15]. In order to give an insight into the overhead introduced by a first-order masking, we implemented two secure adders in C language, using the same compilation options as described in Sect. 4.3. This allowed us to compare, in terms of performance, our secure implementations of ChaCha20 with the one from the OpenSSL library. Running times given in Table 1 are expressed in clock cycles and were computed with the help of debug sessions. Note that these measurements do not take the generation of random numbers into account since this

operation depends a lot on the computing platform. As these countermeasures were implemented in C, they do not ensure the absence of memory accesses within QRs. On the other hand, handling all data in registers during a whole QR may not be possible, since masking also increases memory requirements. Further investigations need to be carried out to determine which algorithms could minimize memory access within QRs and how to securely manage them.

Table 1. Running time in clock cycles to encrypt a 512-bit block using ChaCha20 on an ARM Cortex-M3

	Time	Penalty factor
ChaCha20 unmasked	4 380	1
ChaCha20 with Karroumi <i>et al.</i> SecAdd [22]	121 618	28
ChaCha20 with Coron <i>et al.</i> SecAdd [15]	93 993	22

These practical results point out how difficult it is to effectively secure ARX ciphers’ implementations. However, masking is not the only answer to side-channel attacks and is often combined with *hiding* countermeasures. The principle of hiding is to randomize an algorithm execution by running its operations at different moments in time, during each execution [36, 39]. This can be achieved by randomly inserting dummy operations and *shuffling*. Shuffling intends to randomly change the sequence of operations that can be computed in arbitrary order. In practice, hiding countermeasures increase the number of traces needed to carry out an attack [16, 33].

Regarding ChaCha, operations within a QR cannot be shuffled as they are executed sequentially. On the other hand, each QR can be computed independently from the other, but this is only true for a single round because of switching from column to diagonal rounds. However, there are many ways to implement hiding in practice and further investigations will have to be carried out on the specific case of ChaCha.

8 Conclusions and Further Work

This paper presents side-channel analyses of ChaCha based on leakages related to memory accesses. Our study emphasizes that quantifying the signal available to the attacker at the instruction level could allow to strengthen implementations without much effort.

We compare, from a side-channel point of view, two different software implementations of ChaCha20 on a 32-bit processor. As a result, minimizing memory accesses makes selection functions more complex, to such an extent that they may lead to collisions. We introduce the bricklayer attack to defeat such implementations. Our results show that attacking the reverse QRs (*i.e.* from the keystream) is more efficient than attacking the regular ones (*i.e.* from the input

block). However, we highlight that attacks from the input block are the most pragmatic threats since the knowledge of the keystream is a strong assumption. Finally, we discuss possible countermeasures at several levels and highlight how expensive it is to implement first-order masking for ChaCha20 with practical measurements. Therefore, further work must be undertaken to propose efficient secure implementations of ChaCha.

References

1. chacha20-poly1305@openssh.com: Authenticated encryption mode, May 2016. <http://bvx.su/OpenBSD/usr/bin/ssh/PROTOCOL.chacha20poly1305>
2. iOS 10 Security White Paper. Technical report, Apple Inc., March 2017. https://www.apple.com/business/docs/iOS_Security_Guide.pdf
3. Aumasson, J.-P., Fischer, S., Khazaei, S., Meier, W., Rechberger, C.: New features of Latin dances: analysis of Salsa, ChaCha, and Rumba. Cryptology ePrint Archive, Report 2007/472 (2007). <http://eprint.iacr.org/2007/472>
4. Babbage, S., Borghoff, J., Velichkov, V.: The eSTREAM portfolio in 2012. <http://www.ecrypt.eu.org/ecrypt2/documents/D.SYM.10-v1.pdf>
5. Beaulieu, R., Shors, D., Smith, J., Treatman-Clark, S., Weeks, B., Wingers, L.: SIMON and SPECK: block ciphers for the Internet of Things. Cryptology ePrint Archive, Report 2015/585 (2015). <http://eprint.iacr.org/2015/585>
6. Bernstein, D.J.: Extending the Salsa20 nonce. <https://cr.yp.to/snuffle/xsalsa-20081128.pdf>
7. Bernstein, D.J.: ChaCha, a variant of Salsa20. In: SASC - The State of the Art of Stream Ciphers, pp. 273–278 (2008). <http://cr.yp.to/chacha/chacha-20080128.pdf>
8. Bernstein, D.J.: The Salsa20 family of stream ciphers. In: Robshaw, M., Billet, O. (eds.) *New Stream Cipher Designs*. LNCS, vol. 4986, pp. 84–97. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-68351-3_8
9. Biryukov, A., Dinu, D., Großschädl, J.: Correlation power analysis of lightweight block ciphers: from theory to practice. In: Manulis, M., Sadeghi, A.-R., Schneider, S. (eds.) *ACNS 2016*. LNCS, vol. 9696, pp. 537–557. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-39555-5_29
10. Boura, C., Lvque, S., Vigilant, D.: Side-channel analysis of Grostl and Skein. In: 2012 IEEE Symposium on Security and Privacy Workshops, pp. 16–26, May 2012. <https://www.ieee-security.org/TC/SPW2012/proceedings/4740a016.pdf>
11. Brier, E., Clavier, C., Olivier, F.: Correlation power analysis with a leakage model. In: Joye, M., Quisquater, J.-J. (eds.) *CHES 2004*. LNCS, vol. 3156, pp. 16–29. Springer, Heidelberg (2004). https://doi.org/10.1007/978-3-540-28632-5_2
12. Bursztein, E.: Speeding up and strengthening HTTPS connections for Chrome on Android. Technical report, April 2014. <https://security.googleblog.com/2014/04/speeding-up-and-strengthening-https.html>
13. Callan, R., Zajić, A., Prvulovic, M.: A practical methodology for measuring the side-channel signal available to the attacker for instruction-level events. In: *Proceedings of the 47th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO-47*, pp. 242–254. IEEE Computer Society, Washington, D.C. (2014). <http://dx.doi.org/10.1109/MICRO.2014.39>
14. Choudhuri, A.R., Maitra, S.: Differential cryptanalysis of Salsa and ChaCha - an evaluation with a hybrid model. Cryptology ePrint Archive, Report 2016/377 (2016). <http://eprint.iacr.org/2016/377>

15. Coron, J.-S., Großschädl, J., Tibouchi, M., Vadnala, P.K.: Conversion from arithmetic to boolean masking with logarithmic complexity. In: Leander, G. (ed.) FSE 2015. LNCS, vol. 9054, pp. 130–149. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-48116-5_7
16. Couroussé, D., Barry, T., Robisson, B., Jaillon, P., Potin, O., Lanet, J.-L.: Runtime code polymorphism as a protection against side channel attacks. Cryptology ePrint Archive, Report 2017/699 (2017). <http://eprint.iacr.org/2017/699>
17. Denis, F.: The XChaCha20-Poly1305 construction. https://download.libsodium.org/doc/secret-key_cryptography/xchacha20-poly1305.construction.html
18. Dworkin, M.J.: SP 800-38A 2001 edition. Recommendation for Block Cipher Modes of Operation: Methods and Techniques. Technical report, Gaithersburg, MD, United States (2001)
19. Gandolfi, K., Mourtel, C., Olivier, F.: Electromagnetic analysis: concrete results. In: Koç, Ç.K., Naccache, D., Paar, C. (eds.) CHES 2001. LNCS, vol. 2162, pp. 251–261. Springer, Heidelberg (2001). https://doi.org/10.1007/3-540-44709-1_21
20. Goubin, L.: A sound method for switching between boolean and arithmetic masking. In: Koç, Ç.K., Naccache, D., Paar, C. (eds.) CHES 2001. LNCS, vol. 2162, pp. 3–15. Springer, Heidelberg (2001). https://doi.org/10.1007/3-540-44709-1_2
21. Jungk, B., Bhasin, S.: Don't fall into a trap: physical side-channel analysis of chacha20-poly1305. In: Design, Automation Test in Europe Conference Exhibition (DATE 2017), pp. 1110–1115, March 2017
22. Karroumi, M., Richard, B., Joye, M.: Addition with blinded operands. In: Prouff, E. (ed.) COSADE 2014. LNCS, vol. 8622, pp. 41–55. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10175-0_4
23. Kocher, P., Jaffe, J., Jun, B.: Differential power analysis. In: Wiener, M. (ed.) CRYPTO 1999. LNCS, vol. 1666, pp. 388–397. Springer, Heidelberg (1999). https://doi.org/10.1007/3-540-48405-1_25. <http://dl.acm.org/citation.cfm?id=646764.703989>
24. Kumar, S.V.D., Patranabis, S., Breier, J., Mukhopadhyay, D., Bhasin, S., Chattopadhyay, A., Bakshi, A.: A practical fault attack on ARX-like ciphers with a case study on ChaCha20. In: 2017 Workshop on Fault Diagnosis and Tolerance in Cryptography, FDTCT, Taipei, Taiwan (2017)
25. Langley, A., Chang, W., Mavrogiannopoulos, N., Strombergson, J., Josefsson, S.: ChaCha20-Poly1305 cipher suites for transport layer security (TLS). RFC 7905, RFC Editor, June 2016. <http://tools.ietf.org/rfc/rfc7905.txt>
26. Lemke, K., Schramm, K., Paar, C.: DPA on n -bit sized boolean and arithmetic operations and its application to IDEA, RC6, and the HMAC-construction. In: Joye, M., Quisquater, J.-J. (eds.) CHES 2004. LNCS, vol. 3156, pp. 205–219. Springer, Heidelberg (2004). https://doi.org/10.1007/978-3-540-28632-5_15. <http://www.iacr.org/archive/asiacrypt2007/31560191/31560191.pdf>
27. Luykx, A., Mennink, B., Paterson, K.G.: Analyzing multi-key security degradation. Cryptology ePrint Archive, Report 2017/435 (2017). <http://eprint.iacr.org/2017/435>
28. Maitra, S.: Chosen IV cryptanalysis on reduced round ChaCha and Salsa. Discrete Appl. Math. **208**(C), 88–97 (2016). <http://dx.doi.org/10.1016/j.dam.2016.02.020>
29. Mazumdar, B., Ali, S.S., Sinanoglu, O.: Power analysis attacks on ARX: an application to Salsa20. In: 2015 IEEE 21st International On-line Testing Symposium (IOLTS), pp. 40–43, July 2015

30. McCann, D., Eder, K., Oswald, E.: Characterising and comparing the energy consumption of side channel attack countermeasures and lightweight cryptography on embedded devices. *Cryptology ePrint Archive*, Report 2015/832 (2015). <http://eprint.iacr.org/2015/832>
31. McGrew, D., Bailey, D.: AES-CCM cipher suites for transport layer security (TLS). RFC 6655, RFC Editor, July 2012. <http://tools.ietf.org/rfc/rfc6655.txt>
32. Mozaffari-Kermani, M., Azarderakhsh, R.: Reliable hash trees for post-quantum stateless cryptographic hash-based signatures. In: 2015 IEEE International Symposium on Defect and Fault Tolerance in VLSI and Nanotechnology Systems (DFTS), pp. 103–108, October 2015
33. Patranabis, S., Roy, D.B., Vadnala, P.K., Mukhopadhyay, D., Ghosh, S.: Shuffling across rounds: a lightweight strategy to counter side-channel attacks. In: 2016 IEEE 34th International Conference on Computer Design (ICCD), pp. 440–443, October 2016
34. Quisquater, J.-J., Samyde, D.: ElectroMagnetic Analysis (EMA): measures and counter-measures for smart cards. In: Attali, I., Jensen, T. (eds.) *E-smart 2001*. LNCS, vol. 2140, pp. 200–210. Springer, Heidelberg (2001). https://doi.org/10.1007/3-540-45418-7_17
35. Rescorla, E.: The transport layer security (TLS) protocol version 1.3. Internet-Draft draft-ietf-tls-tls13-21, Internet Engineering Task Force, July 2017. <https://tswg.github.io/tls13-spec/draft-ietf-tls-tls13.html>, work in Progress
36. Rivain, M., Prouff, E., Doget, J.: Higher-order masking and shuffling for software implementations of block ciphers. *Cryptology ePrint Archive*, Report 2009/420 (2009). <http://eprint.iacr.org/2009/420>
37. Salowey, J., Choudhury, A., McGrew, D.: AES Galois Counter Mode (GCM) cipher suites for TLS. RFC 5288, RFC Editor, August 2008. <http://www.rfc-editor.org/rfc/rfc5288.txt>
38. Shi, Z., Zhang, B., Feng, D., Wu, W.: Improved key recovery attacks on reduced-round Salsa20 and ChaCha. In: Kwon, T., Lee, M.-K., Kwon, D. (eds.) *ICISC 2012*. LNCS, vol. 7839, pp. 337–351. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-37682-5_24
39. Veyrat-Charvillon, N., Medwed, M., Kerckhof, S., Standaert, F.-X.: Shuffling against side-channel attacks: a comprehensive study with cautionary note. In: Wang, X., Sako, K. (eds.) *ASIACRYPT 2012*. LNCS, vol. 7658, pp. 740–757. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-34961-4_44
40. Yadav, P., Gupta, I., Murthy, S.K.: Study and analysis of eSTREAM cipher Salsa and ChaCha. In: 2016 IEEE International Conference on Engineering and Technology (ICETECH), pp. 90–94, March 2016
41. Zohner, M., Kasper, M., Stöttinger, M.: Butterfly-attack on Skein’s modular addition. In: Schindler, W., Huss, S.A. (eds.) *COSADE 2012*. LNCS, vol. 7275, pp. 215–230. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-29912-4_16