

# Flexible Exponentiation with Resistance to Side Channel Attacks

Camille Vuillaume and Katsuyuki Okeya

Hitachi, Ltd., Systems Development Laboratory, Kawasaki, Japan  
{camille, ka-okeya}@sdl.hitachi.co.jp

**Abstract.** We present a countermeasure for protecting modular exponentiations against side-channel attacks such as power, timing or cache analysis. Our countermeasure is well-suited for tamper-resistant implementations of RSA or DSA, without significant penalty in terms of speed compared to commonly implemented methods. Thanks to its high efficiency and flexibility, our method can be implemented on various platforms, from smartcards with low-end processors to high-performance servers.

**Keywords:** RSA, side channel attacks, fractional width, simple power analysis.

## 1 Introduction

With the rise of electronic communications, and in particular, electronic commerce, public-key cryptography has become an essential component in our daily life. The de-facto standard for public-key encryption and digital signatures is RSA, and with the development of miniaturization, RSA is now implemented not only on high-performance servers, but also on various mobile devices such as smartcards or mobile phones.

It is believed that with a bitlength of 1024 bits, RSA is secure for middle-term applications, and with 2048 bits, for long-term applications. However, protecting against mathematical attacks is not sufficient in the real world. Indeed, it has been shown that practical implementations of cryptosystems often suffer from critical information leakage through side-channels: timings [6], power consumption [5] or cache usage [13], for instance. Such side-channel attacks are no theoretical works that researchers secretly run in laboratories with expensive hardware, but practical threats to virtually any application where secrecy matters.

On the one hand, there are numerous countermeasures for defeating side-channel attacks on elliptic curve cryptosystems [4, 7, 9], but on the other hand, there are few of them for RSA [15, 16]. At first sight, it seems that the keylength of RSA is so long that even when side-channel information partially reveals the secret, exhaustive searches often remain ineffective. However, when a sufficiently large part of the secret key is known, RSA can be broken [3]. Thus, despite long secret keys, RSA also needs a decent protection against side-channel attacks [13].

Unfortunately, many countermeasures that have been developed for elliptic curves cannot be transposed to the case of RSA: on elliptic curves, signed representations of the exponents are heavily used because they yield faster exponentiation algorithms when inversions are cheap. But unlike elliptic curves, RSA does not benefit from cheap inversions, and therefore, unsigned representations are the only option [16]. In other words, countermeasures which are efficient on elliptic curves are generally not practical at all in the case of RSA.

Our contribution is as follows: we show how to transform a countermeasure based on a signed representation [7] into a countermeasure based on an unsigned representation. Then, we improve the flexibility of our countermeasure with a fractional width technique [10]: with the improved method, the size of the pre-computed table can be freely chosen. Our countermeasure is not only highly flexible, and therefore well-suited for a wide range of platforms, from constrained environments such as smartcards to high-performance servers, but in addition, can compete with the commonly implemented exponentiation techniques. In practical situations, our method is only about 5% slower than commonly used methods. And most importantly, our method thwarts several types of attacks, power analysis and cache attacks in particular. Finally, we refine attacks against fractional width techniques and introduce a new tool to evaluate the quality of our countermeasure in the sense of resistance to side channel attacks. We demonstrate that in the SPA model, our countermeasure can prevent information leakage.

## 2 Side Channel Attacks

Over the past few years, theoretical attacks against cryptosystems such as RSA or elliptic curves have shown little improvements. On the other hand, attacks based on physical information leakage, also known as side channel attacks, have revolutionized the definition of tamper-resistance.

### 2.1 Methodology of Side Channel Attacks

Side channel attacks take advantage of the correlation between secret values and physical emanations such as timings [6], power consumption [5] or even cache usage [13].

*Power Analysis.* On smartcards, which do not have any embedded power supply, the most powerful approach is probably to measure the power consumption of the device supplied from the outside [5]. One can classify power analysis attacks into two main classes: simple power analysis (SPA) and differential power analysis (DPA). The approach of SPA is to identify regions of a power trace which directly depend on the secret key. It is common for cryptographic algorithms to have conditional branches depending on the value of some secret bit, and typically, those branches are targets of choice for SPA. For example, in the frame of an RSA exponentiation computed with the binary method, the binary representation of the (secret) exponent is scanned; for the bit-value zero, a square is computed,

whereas a square and a multiplication are calculated when the bit-value is one. Thus, it is easy to see that the knowledge of the operation sequence (square or square-multiplication) is equivalent to the knowledge of the secret exponent. In other words, by distinguishing squares from multiplications in power traces, one can reveal the secret exponent. A similar problem also exists in the case of elliptic curves, where the ability of distinguishing the two types of elliptic operations, namely point doublings and point additions, also leads to the secret scalar. DPA is more sophisticated: the idea is to guess the value of the secret bit-by-bit, and try to confirm or infirm the guess for each bit thanks to statistical analysis of several power traces.

*Cache Attacks.* Recently, cache attacks on public key cryptosystems have been investigated, and proof-of-concept attacks based on cache analysis confirmed that the cache should not be neglected as source of information leakage [13]. On computers, power analysis is extremely invasive, and although possible from a theoretical point of view, difficult to set up in practice. On the contrary, cache attacks are practical, because cryptographic algorithms leave characteristic footprints in the cache memory and other processes can spy on the cache. Using such techniques, it has been shown that on computers equipped with the hyperthreading technology, using one single RSA exponentiation, cache observation reveals the secret key [13].

## 2.2 Countermeasures

Countermeasures against side-channel attacks have been proposed on various cryptosystems and with various strategies. In particular, alternative representations of the secret are often used as countermeasures, on elliptic curves and on RSA.

*Strategies and Requirements.* It is not that difficult to protect cryptosystems against DPA. By definition, DPA requires that the same secret is used to perform several cryptographic operations with each time a different input value: decrypting or signing several messages, for instance. Signature schemes such as DSA and EC-DSA use a new random ephemeral as exponent for each new signature, and as a consequence, are naturally immune to DPA. In the case of RSA, a well-known technique to defeat DPA is to blind the secret exponent, that is, to add a random multiple of the group order to the secret. Because  $c^{d+r\phi(n)} = c^d \bmod n$ , blinding does not change the result of the exponentiation, but in the same time, changes the exponent itself. As a consequence, the side-channel information that arises from multiple power traces is not correlated. On the contrary, despite the relative simplicity of the idea behind SPA, it is not easy to design secure and efficient SPA countermeasures. However, SPA-resistance is always necessary, and is a prerequisite to DPA resistance. For instance, if only one random ephemeral exponent of DSA or EC-DSA is revealed, the secret key of the signature scheme can be easily inferred. Similarly, from the point of view of the attacker, a blinded RSA exponent  $d + r\phi(n)$  is as good as the secret itself.

In the following, we explain how the *representation* of the secret exponent can be an effective countermeasure against side channel attacks, and especially against SPA. We call representation a function  $\mathcal{R} : D^k \rightarrow \{0, 1, \dots, 2^\mathcal{L} - 1\}$ , where  $D$  is the *digit set* of the representation. For instance, the binary representation  $\mathcal{R}_b(b_{\mathcal{L}-1} \dots b_0) \in \{0, 1\}^\mathcal{L} \mapsto \sum_{i=0}^{\mathcal{L}-1} b_i 2^i \in \{0, 1, \dots, 2^\mathcal{L} - 1\}$  utilizes base-2 digits (bits) to represent integers. Rather than  $\{0, 1\}$ , larger digit sets are sometimes preferred: in that case, pre-computations allow a faster execution of the cryptographic primitive. Finally, signed representations such as  $\{0, \pm 1\}$  are advantageous when the computation of inverses is easy. In particular, this is attractive on elliptic curves where  $-P$  can be computed from  $P$  for almost free. But the case of RSA is different: inversions are very expensive, and signed representations do not yield any interest in practice for RSA. Some representations are also valuable countermeasures against side channel attacks. Indeed, by changing the *representation* of the secret without changing its value, one can gain control on the operation sequence, and ultimately, on side-channel information leakage.

*Möller's Countermeasure on Elliptic Curves.* On elliptic curve cryptosystems, it is common to use representations with a large digit set and to pre-compute some small multiples of the base point in order to speed up the scalar multiplication. In particular, thanks to window methods with digit set  $\{0, 1, 2, \dots, 2^w - 1\}$ , one can reduce the computational cost of the scalar multiplication given that small multiples of the base point  $P$  are pre-computed:  $2P, 3P, \dots, (2^w - 1)P$ . Although standard window methods aim at greater efficiency only, they can also be enhanced to SPA-resistant scalar multiplication schemes, where the secret is recoded with a fixed pattern, using a signed representation with digit set  $\{-2^w, 1, 2, \dots, 2^w - 1\}$ , where 0 is absent [7]. More precisely, an additional conversion step is applied in each window in order to remove zero digits:

1. replace the digit 0 with  $-2^w$ , and add a carry of +1 to the next window when scanning the scalar from right to left,
2. replace the digit  $2^w$  with  $-2^w$ , and add a carry of +2 to the next window,
3. replace the digit  $2^w + 1$  with 1 and add a carry of +1 to the next window,
4. otherwise leave the digit as it is.

Thanks to this conversion technique, the scalar can be recoded with a fixed pattern: nonzero digits are always followed by exactly  $w - 1$  zero digits. The advantage of this approach is that the operation pattern becomes regular as well: point additions are always followed by exactly  $w$  point doublings, which makes SPA impractical. Because the point  $-2^w P$  must be readily available, the scheme requires that computing inversions in the group (in that case, elliptic point negative) is easy. Unfortunately, this is not the case for RSA.

*Countermeasures on RSA.* An obvious way to thwart SPA on RSA and other cryptosystems is to insert dummy operations in such a way that the operation sequence does not depend on the secret anymore. But this technique suffers from

severe drawbacks: not only the computational cost is considerably increased, but the use of dummy operations is generally not recommended due to safe-error attacks, which can take advantage of such dummy operations [17]. Side-channel atomicity is a more sophisticated countermeasure, where side-channel information consists of the repetition of an atomic side-channel block [2]. However, on RSA, side-channel atomicity requires multiplications and squares to be indistinguishable. This is often not the case on high-speed implementations of RSA in software, where squares are implemented with a distinct procedure. Furthermore, side-channel atomicity does not address cache attacks.

In addition to that, there are some countermeasures for RSA which are based on randomized representations of the secret [15,16]; their aim is to protect against SPA and in the same time to improve resistance to differential attacks. The principle of the MIST exponentiation algorithm [15] is to randomly change the basis (in other words, the digit set  $D$ ) during the recoding. Rather than a pre-computed table, MIST utilizes efficient addition chains and is faster than the binary method. Following a different approach, Yen et al. utilize a large pre-computed table and randomize the representation of the secret [16]. Their exponentiation scheme is computed from left to right, without need for any inversion. With 14 pre-computed values, the efficiency of the countermeasure is about the same as the window method with  $w = 2$ . However, both of these methods suffer from a considerable overhead compared to high-speed techniques. For instance, compared to the sliding window method with  $w = 5$ , which is used in OpenSSL, these countermeasures yield more than 30% performance drop.

### 3 SPA-Resistant Unsigned Recoding Techniques

On elliptic curves, numerous recoding techniques have been proposed for means of defeating side-channel attacks. These representations are often based on signed digit sets. Unfortunately, this approach is not valid for RSA where computing inverses is too costly; hence our motivation to construct secure unsigned representations.

#### 3.1 SPA-Resistant Unsigned Integral Width

Our approach is to extend Möller's recoding to the unsigned case. To obtain the unsigned digit set  $\{1, 2, \dots, 2^w\}$ , the key idea of our method is to use *negative* carries rather than positive carries:

1. replace the digit 0 with  $2^w$ , and add a carry of  $-1$  to the next window when scanning the scalar from right to left,
2. replace the digit  $-1$  with  $2^w - 1$ , and add a carry of  $-1$  to the next window,
3. otherwise leave the digit as it is.

On the one hand, it is easy to see that Möller's algorithm terminates, and that if the original bitlength of the scalar was  $\mathcal{L}$ , the recoded scalar has at most

$\mathcal{L} + 1$  digits. On the other hand, in the case of the above rules for generating an unsigned representation with a fixed pattern, the situation is different: there is no guarantee that the algorithm will terminate because a carry can propagate indefinitely. To ensure a correct termination, we treat the case of the most significant bit separately: if a carry remains at the end of the recoding, we use the most significant bit to neutralize it, and reduce the length of the exponent. If instead of  $d$ , a blinded exponent  $d + r\phi(n)$  with random  $r$  is used as DPA countermeasure, this approach is safe<sup>1</sup>. But if not, there is a direct information leakage, because the length of the recoded exponent depends on the value of some secret bits. To remove this leakage, we extend the bitlength of the exponent by 2, and fix the value of the 2 most significant bits  $d_{\mathcal{L}+1} = 1$  and  $d_{\mathcal{L}} = 0$ . This is always possible because  $c^{d+\phi(n)} = c^d \bmod n$ : in other words, adding  $\phi(n)$  to the exponent does not change the value of the exponentiation modulo  $n$ . By repeatedly adding  $\phi(n)$ , one can always set  $d_{\mathcal{L}+1}$  to 1 and  $d_{\mathcal{L}}$  to 0 (because  $d < \phi(n) < n < 2^{\mathcal{L}}$ ). If the value of  $d$  and  $n$  are fixed (which is typically the case when  $d$  is a secret key), this calculation can be performed once for all at the key generation stage. Now, since  $d_{\mathcal{L}} = 0$ , independently from the value of the previous bits, the corresponding recoded digit is  $u_{\mathcal{L}} \neq 0$  and a carry  $\gamma$  is generated. Finally,  $u_{\mathcal{L}+1} = d_{\mathcal{L}+1} - 1 = 0$ , therefore the length of the recoded expansion is *always* reduced by one.

---

**Algorithm 1.** Conversion to unsigned integral SPA-resistant representation

---

INPUT:  $\mathcal{L} + 2$ -bit exponent  $d = (10d_{\mathcal{L}-1} \dots d_0)_2$ , width  $w$ ;

OUTPUT: Recoded exponent  $(u_{\mathcal{L}} \dots u_0)$ ;

---

1.  $i \leftarrow 0$ ;  $\gamma \leftarrow 0$ ;
  2. **while**  $i \leq \mathcal{L} - w$  **do**
    - (a)  $u_i \leftarrow (d_{i+w-1} \dots d_i)_2 - \gamma$ ;
    - (b) **if**  $u_i \leq 0$  **then**  $u_i \leftarrow u_i + 2^w$ ;  $\gamma \leftarrow 1$ ; **else**  $\gamma \leftarrow 0$ ;
    - (c)  $u_{i+1} \leftarrow 0, \dots, u_{i+w-1} \leftarrow 0$ ;  $i \leftarrow i + w$ ;
  3. **if**  $i < \mathcal{L}$  **then**
    - (a)  $u_i \leftarrow (d_{\mathcal{L}-1} \dots d_i)_2 - \gamma$ ;
    - (b) **if**  $u_i \leq 0$  **then**  $u_i \leftarrow u_i + 2^{\mathcal{L}-i}$ ;  $\gamma \leftarrow 1$ ; **else**  $\gamma \leftarrow 0$ ;
    - (c)  $u_{i+1} \leftarrow 0, \dots, u_{\mathcal{L}-1} \leftarrow 0$ ;
  4.  $u_{\mathcal{L}} \leftarrow 2 - \gamma$ ; **return**  $(u_{\mathcal{L}} \dots u_0)$ ;
- 

*Exponentiation With Integral Width Recoding.* Algorithm 2 computes  $g^d$  with an unsigned SPA-resistant recoding of  $d$ . In fact, during the exponentiation, the operation pattern is fixed:  $w$  squares and 1 multiplication. Note that pre-computations take advantage of faster squares.

---

<sup>1</sup> Blinding alone is not sufficient, because it protects only against DPA, not SPA, and should be used with an additional SPA countermeasure such as our recoding.

---

**Algorithm 2.** Integral width pre-computations and exponentiation

---

INPUT:  $\mathcal{L}$  + 2-bit exponent  $d = (10d_{\mathcal{L}-1} \dots d_0)_2$ , width  $w$ , basis  $g$ ;OUTPUT:  $c = g^d$ ;

---

1.  $g[1] \leftarrow g$ ;
  2. **for**  $i$  **from** 2 **to**  $2^w$  **step** 2 **do**  $g[i] \leftarrow g[i/2]^2$ ;  $g[i+1] \leftarrow g[i] * g$ ;
  3. recode  $d$  with Algorithm 1;  $i \leftarrow \mathcal{L} - 1$ ;  $c \leftarrow g[u_{\mathcal{L}}]$ ;
  4. **while**  $i \geq 0$  **do**
    - (a)  $c \leftarrow c^2$ ;
    - (b) **if**  $u_i > 0$  **then**  $c \leftarrow c * g[u_i]$ ;
    - (c)  $i \leftarrow i - 1$ ;
  5. **return**  $c$ ;
- 

### 3.2 SPA-Resistant Unsigned Fractional Width

A disadvantage of the previous method, and more generally, of table-based exponentiations, is that there are only limited choices for the table size. Since the table size and the cost of pre-computations grow exponentially with the width  $w$ , number of bits which are scanned simultaneously, large values of  $w$  become quickly impractical. However, it would be useful to be able to select *any* table size. Fractional width recodings make this possible, thanks to a degenerated width- $w$  pre-computed table where some values are missing [8]. In addition, SPA-resistant fractional width methods exist in the case of elliptic curves, where it is advantageous to use signed representations [10]. We show that our unsigned (integral width) SPA-resistant recoding technique (Algorithm 1) can also be enhanced to an unsigned fractional width recoding.

*Unsigned Fractional Width Recoding.* The idea is the same as in the original signed SPA-resistant fractional width recoding: compute simultaneously the digits  $x$  and  $y$ , where  $x$  correspond to the width  $w$  and  $y$  to the width  $w - 1$ , and choose  $x$  or  $y$  depending on some criteria. The knowledge of the choice of  $x$  or  $y$  must not help attackers to gather information about the secret. To fulfill this requirement, the set of pre-computed values is randomized. More precisely, when the pre-computed table has  $k$  entries, define  $w = \lceil \log_2(k) \rceil$ . Then, the  $2^{w-1}$  elements  $\{c, c^2, \dots, c^{2^{w-1}}\}$  are always pre-computed, but  $k - 2^{w-1}$  additional pre-computed elements are randomly chosen in the set  $\{c^{2^{w-1}+1}, c^{2^{w-1}+2}, \dots, c^{2^w}\}$ . In the following, we call  $B$  the set of the exponents of the chosen pre-computed elements.

The core idea in Algorithm 3 is that the choice of recoding a sequence of bits with the width  $w$  or  $w - 1$  looks random to the attacker. The principle of the recoding is the same as in the signed fractional width algorithm [10]. Recall that  $B$  is the set of exponents of the pre-computed values, define the width  $w = \lceil \log_2(k) \rceil$  and the probability  $p = k/2^{w-1} - 1$ . Then, for  $x$  computed with width  $w$  and  $y$  with width  $w - 1$ , we apply the following rules:

---

**Algorithm 3.** Conversion to unsigned fractional SPA-resistant representation

---

INPUT:  $\mathcal{L}$  + 2-bit exponent  $d = (10d_{\mathcal{L}-1} \dots d_0)_2$ , table size  $k$ , index set  $B$ ;

OUTPUT: Recoded exponent  $(u_{\mathcal{L}} \dots u_0)$ ;

---

1.  $i \leftarrow 0$ ;  $\gamma \leftarrow 0$ ;  $w \leftarrow \lceil \log_2(k) \rceil$ ;
  2. **while**  $i \leq \mathcal{L} - w$  **do**
    - (a)  $x \leftarrow (d_{i+w-1} \dots d_i)_2 - \gamma$ ;
    - (b)  $y \leftarrow (d_{i+w-2} \dots d_i)_2 - \gamma$ ;
    - (c) **if**  $x \leq 0$  **then**  $x \leftarrow x + 2^w$ ;  $\gamma_x \leftarrow 1$ ; **else**  $\gamma_x \leftarrow 0$ ;
    - (d) **if**  $y \leq 0$  **then**  $y \leftarrow y + 2^{w-1}$ ;  $\gamma_y \leftarrow 1$ ; **else**  $\gamma_y \leftarrow 0$ ;
    - (e) **if**  $x \leq 2^{w-1}$  **then**
      - i. *rnd*  $\leftarrow$  generate  $w - 1$  random bits;
      - ii. **if**  $\text{rnd} < k - 2^{w-1}$  **then**  $u_i \leftarrow x$ ;  $\gamma \leftarrow \gamma_x$ ;  $r \leftarrow w$ ;
      - iii. **else**  $u_i \leftarrow y$ ;  $\gamma \leftarrow \gamma_y$ ;  $r \leftarrow w - 1$ ;
    - (f) **else if**  $x \in B$  **then**  $u_i \leftarrow x$ ;  $\gamma \leftarrow \gamma_x$ ;  $r \leftarrow w$ ;
    - (g) **else**  $u_i \leftarrow y$ ;  $\gamma \leftarrow \gamma_y$ ;  $r \leftarrow w - 1$ ;
    - (h)  $u_{i+1} \leftarrow 0, \dots, u_{i+r-1} \leftarrow 0$ ;  $i \leftarrow i + r$ ;
  3. **if**  $i < \mathcal{L}$  **then**
    - (a)  $u_i \leftarrow (d_{\mathcal{L}-1} \dots d_i)_2 - \gamma$ ;
    - (b) **if**  $u_i \leq 0$  **then**  $u_i \leftarrow u_i + 2^{\mathcal{L}-i}$ ;  $\gamma \leftarrow 1$ ; **else**  $\gamma \leftarrow 0$ ;
    - (c)  $u_{i+1} \leftarrow 0, \dots, u_{\mathcal{L}-1} \leftarrow 0$ ;
  4.  $u_{\mathcal{L}} \leftarrow 2 - \gamma$ ; **return**  $(u_{\mathcal{L}} \dots u_0)$ ;
- 

1. if  $x \leq 2^{w-1}$  then choose  $x$  with probability  $p$  or  $y$  with probability  $1 - p$ ,
2. else if  $x \in B$  (in other words,  $g^x$  is pre-computed), choose  $x$  (this occurs with probability  $p$ ),
3. else choose  $y$  (in that case,  $g^x$  is not pre-computed, this happens with probability  $1 - p$ ).

Therefore, for randomly chosen exponents, the width  $w$  is chosen with probability  $p$  and  $w - 1$  with probability  $1 - p$ . Additionally, since the set  $B$  is randomized for each new recoding, the two patterns can actually appear for the same sequence of bits.

*Exponentiation With Fractional Width Recoding.* The technique for the exponentiation with a fractional SPA-resistant width is almost the same as in the integral case, with the exception of the pre-computation step. Indeed, since the pre-computed table is de-generated, the values in the upper half part of the table are computed with a special procedure. Since  $g^i = g^{i-2^{w-1}} * g^{2^{w-1}}$ , the pre-computations use the value  $g^{2^{w-1}}$  which is already available in order to compute the upper half part of the table. Note that in Algorithm 4, for the sake of simplicity, the pre-computed table is indexed with the exponent of the pre-computed values. In reality, a look-up table should be used in order to save memory: such table could be indexed with the exponents of pre-computed values, and would additionally store a pointer to the actual location of the pre-computed value.

**Algorithm 4.** Fractional width pre-computations and exponentiation

INPUT:  $\mathcal{L} + 2$ -bit exponent  $d = (10d_{\mathcal{L}-1} \dots d_0)_2$ , table length  $k$ , basis  $g$ ;  
 OUTPUT:  $c = g^d$ ;

1.  $g[1] \leftarrow g$ ;
2. **for**  $i$  **from** 2 **to**  $2^{w-1}$  **step** 2 **do**  $g[i] \leftarrow g[i/2]^2$ ;  $g[i+1] \leftarrow g[i] * g$ ;
3. **for all**  $i > 2^{w-1}$ ,  $i \in B$  **do**  $g[i] \leftarrow g[i - 2^{w-1}] * g[2^{w-1}]$ ;
4. recode  $d$  with Algorithm 3;  $i \leftarrow \mathcal{L} - 1$ ;  $c \leftarrow g[u_{\mathcal{L}}]$ ;
5. **while**  $i \geq 0$  **do**
  - (a)  $c \leftarrow c^2$ ;
  - (b) **if**  $u_i > 0$  **then**  $c \leftarrow c * g[u_i]$ ;
  - (c)  $i \leftarrow i - 1$ ;
6. **return**  $c$ ;

**3.3 Efficiency of the Fractional Width Exponentiation**

Next, we describe the advantages of our technique in terms of speed and memory, and compare its performances with that of commonly used exponentiation methods. For a table size  $k$  and a bitlength  $n$ ,  $k$  elements of  $n$  bits (including the basis of the exponentiation  $g$  are pre-computed and stored in RAM. In contrary to other methods, which only allows 1, 2, 4,  $\dots$ ,  $2^w, \dots$  as table size, the fractional window method is much more flexible: any table size can be chosen. This is not only an advantage to fully occupy the available memory on constrained environments, but also means that the optimal table size  $k$  can be chosen on large-memory profiles: in that case, the fractional width method yields an exponentiation method which is faster than integral width techniques.

*Efficiency.* The cost of the pre-computations of the unsigned SPA-resistant technique is as follows: for the  $2^{w-1}$  pre-computed values in the lower half table,  $2^{w-2}$  squares and  $2^{w-2} - 1$  multiplications, and for the upper half table,  $k - 2^{w-1}$  multiplications. Recall that the upper width  $w$  is defined as  $w = \lceil \log_2(k) \rceil$ , and that probability of choosing the width  $w$  rather than  $w - 1$  is  $p = \frac{k}{2^{w-1}} - 1$ . Multiplications occur only when there is a nonzero digit in the representation, therefore, for an  $n$ -bit exponent, there are on average  $n/(w - 1 + p)$  multiplications. Then, the memory and average computational cost of exponentiations based on the unsigned fractional representations are as follows:

$$\begin{cases} \mathcal{M}_F = k \cdot \mathcal{L} \text{ bits} \\ \mathcal{C}_F = (2^{w-2} + \mathcal{L}) \cdot S + \left( k - 2^{w-2} - 1 + \frac{\mathcal{L}}{w-1+p} \right) \cdot M, \end{cases} \quad (1)$$

where  $S$  and  $M$  stand for the cost of squares and multiplications, respectively,  $w = \lceil \log_2(k) \rceil$  and  $p = \frac{k}{2^{w-1}} - 1$ .

*Comparison with the Sliding Window.* The sliding window method is often utilized for practical implementations of exponentiations; for instance, OpenSSL

**Table 1.** Memory and cost of several exponentiation methods

<b>512-bit exponentiation</b>	Memory	Speed
Binary method	0 bytes	2.563 ms
Sliding window, $w = 5$	1,024 bytes	2.026 ms
Our technique, $k = 16$	1,024 bytes	2.173 ms
Our technique, $k = 33$	2,112 bytes	2.137 ms
<b>1024-bit exponentiation</b>	Memory	Speed
Binary method	0 bytes	15.05 ms
Sliding window, $w = 6$	4,096 bytes	11.49 ms
Our technique, $k = 32$	4,096 bytes	12.09 ms
Our technique, $k = 53$	6,784 bytes	12.05 ms

uses the sliding window with 16 pre-computed values  $\{g, g^3, g^5, \dots, g^{31}\}$  [12]. The cost of pre-computations for the sliding window method with width  $w$ , that is,  $2^{w-1}$  pre-computed values, is one square and  $2^{w-1} - 1$  multiplications. The idea of the sliding window is to consider odd pre-computed values only, reducing the size of the table and the number of multiplications.

$$\begin{cases} \mathcal{M}_{SW} = 2^{w-1} \cdot \mathcal{L} \text{ bits} \\ \mathcal{C}_{SW} = (\mathcal{L} + 1) \cdot S + \left(2^{w-1} - 1 + \frac{\mathcal{L}}{w+1}\right) \cdot M, \end{cases} \quad (2)$$

For the same memory, it is clear that the sliding window is faster than the SPA-resistant fractional window. However, thanks to its higher flexibility, the fractional window has a “better” optimal table size  $k$ . We implemented both techniques with the NTL library [14] and compared the algorithms running with their optimal parameters on our platform; Table 1 summarizes our implementation results. It comes out that, although the fractional width is slightly slower than the sliding window, the performance drop is very small in practice: only 5% for 512-bit and 1024-bit exponentiations and when the algorithms run in their optimal settings. For the same memory consumption, our method is 7% slower than the sliding window 512-bit exponentiation, and 5% than the 1024-bit exponentiation.

## 4 Security Analysis of the Unsigned Fractional Width

We now analyze the security of the unsigned SPA-resistant fractional width technique in the sense of SPA (and related attacks). We particularly study the distributions of digits in the representation, and introduce a tool to measure the quality of the countermeasure.

### 4.1 Non-uniform Digit Distribution

It has been shown that the signed fractional width recoding leaks some information about the secret [11]. The reason for this (relatively small) information leakage is the *degenerated* width  $w$  pre-computed table: the fact that some

pre-computed values are missing in the table can be used to speed up attacks. Obviously, since it uses exactly the same principle for its recoding (but with a different digit set), the unsigned fractional width algorithm (Algorithm 3) can be targeted by this kind of attack as well. In the following, we refine the technique described in [11] by considering the influence of the degenerated table not only on blocks of bits recoded with the width  $w$  but also  $w - 1$ .

*Non-uniformity when  $w$  is selected.* The first point is that when the width  $w$  is selected rather than  $w - 1$ , some digits do not appear in the recoding [11]. Indeed, to construct a length- $k$  pre-computed table, one takes a length- $2^w$  table and remove some values at random from its upper half part. Therefore, the exponents of the removed pre-computed values are absent from the representation. In other words, each time the width  $w$  is selected, there are only  $k$  possible digits instead of  $2^w$ .

For instance, consider  $k = 3$ ; from the length-4 table  $\{g^1, g^2, g^3, g^4\}$ , one value is chosen among  $g^3$  and  $g^4$ , and removed from the table. In our example, we remove  $g^3$ . Then, our representation admits the following digit set:  $B = \{1, 2, 4\}$ .

*Non-uniformity when  $w - 1$  is selected.* The second point is that the distribution of digits recoded with the width  $w - 1$  is not uniform. This problem occurs because when a missing digit is scanned in the secret (case where  $x > 2^{w-1}$  and  $x \notin B$ ), the width  $w - 1$  is selected. But this event occurs only for some digits  $x$ , and as a consequence there are also constraints on the corresponding digits  $y$ . More precisely,  $x = y + 2^{w-1}d_{i+w-1}$ , and for a given value of  $x \notin B$  there is only one possible value for  $y$ . Since there are only  $2^w - k$  for  $x \notin B$ , there are also  $2^w - k$  possible values for  $y$  (instead of  $2^{w-1}$ ). But when  $x \leq 2^{w-1}$  and the width  $w - 1$  is (randomly) chosen, the distributions of digits is uniform, therefore the case where  $x > 2^{w-1}$  and  $x \notin B$  introduces non-uniformity in the distribution of digits recoded with the width  $w - 1$ .

Assume for instance a table size  $k = 5$  and digit set  $B = \{1, 2, 3, 4, 5\}$ . Suppose that during the recoding,  $x = 6$  is computed. Now, there are two possibilities: either  $(d_{i+2}d_{i+1}d_i)_2 = 6$  and the carry is zero ( $\gamma = 0$ ), or  $(d_{i+2}d_{i+1}d_i)_2 = 7$  and the carry is not zero ( $\gamma = 1$ ). In both cases,  $(d_{i+1}d_i)_2 - \gamma = y = 2$ . But when  $x \leq 4$  and  $w - 1 = 2$  is (randomly) selected, the distribution of the four digits  $\{1, 2, 3, 4\}$  is uniform. Consequently, the digit 2 has a higher probability than 1, 3 and 4 when the width  $w - 1 = 2$  is selected. Note that although we concentrate on the case of the unsigned SPA-resistant fractional window in this paper, the signed technique shares similar properties.

## 4.2 Side-Channel Information and Entropy

Although the non-uniformity of the digit distribution is intuitively a problem, even when the bias is known, it is difficult to evaluate how serious the threat is. To explicitly evaluate SPA-resistance of the fractional width techniques, the calculation of the entropy of such representations is helpful.

Consider the following challenge: finding the secret integer  $d \in \{0, 1, \dots, 2^\mathcal{L} - 1\}$ . When performing side-channel analysis, attackers might be able to reveal a bias in the distribution of the candidates in  $\{0, 1, \dots, 2^\mathcal{L} - 1\}$ . The knowledge of this bias can help them find the secret integer by trying out the most probable candidates first.

**Definition 1 (Side-Channel Information Entropy).** *Let  $p$  be a probability function  $p : \{0, 1, \dots, 2^\mathcal{L} - 1\} \rightarrow 0..1$  with  $\sum_{\delta \in \{0, 1, \dots, 2^\mathcal{L} - 1\}} p(\delta) = 1$ , representing a bias in the search space  $\{0, 1, \dots, 2^\mathcal{L} - 1\}$  obtained by side-channel analysis. We call side-channel information entropy the term:*

$$S = - \sum_{\delta \in \{0, 1, \dots, 2^\mathcal{L} - 1\}, p(\delta) \neq 0} p(\delta) \log_2(p(\delta)). \quad (3)$$

Assume for example that an attacker tries to find a secret  $d \in \{0, 1, 2, 3\}$ . Without the help of side-channel analysis, all candidates in the search space are equiprobable:  $p(0) = p(1) = p(2) = p(3) = 1/4$  and the entropy is  $S = 2$  bits. But imagine now that the attacker identified a feature in side-channel information which is more likely to occur for  $\delta = 3$  than for the other candidates. For instance,  $p(3) = 1/2$ , and  $p(0) = p(1) = p(2) = 1/6$ . In that case, the entropy is reduced to  $S = 1.79$  bits.

More generally, when all values for  $d$  are equiprobable even when SPA information is available, the entropy reaches its maxima, namely  $S = \mathcal{L}$  bits. In that case, the representation is a perfect SPA countermeasure. When the values are not uniformly distributed (some values are forbidden or simply less probable),  $S < \mathcal{L}$  bits. One interpretation of the entropy is the number of bits that remain secure (unknown to the attacker) despite side-channel information leakage. Note that the entropy does not necessarily represent the running time of the fastest attacks; in particular, on RSA, factoring the modulus is much faster than running an exhaustive search, even when the search space has been reduced by SPA. However, if too many bits of the secret key are revealed, factoring becomes easy [3]. But if the entropy is large enough so that there exist no practical attack on  $S$  bits, the countermeasure is secure against SPA. Thus, the entropy term  $S$  represents the quality of an SPA countermeasure.

Note that in the case of fractional window methods, the probabilistic recoding process is divided into two groups of events. First, the selection of the pre-computed table; in this case, we assume that the  $\Omega_B$  possible choices for the pre-computed table are *equiprobable* and *indistinguishable* by SPA. And second, the recoding process itself, where some choices between the widths  $w$  and  $w - 1$  are random. We call  $S_w$  and  $S_{w-1}$  the entropy of  $w$ -blocks (that is, a block of bits recoded with width  $w$ ) and  $w - 1$ -blocks. Then, we abuse notations and call “entropy of the fractional width recoding” the term  $S = \log_2 |\Omega_B| + \mathcal{L}_w S_w + \mathcal{L}_{w-1} S_{w-1}$ , where  $\mathcal{L}_w$  and  $\mathcal{L}_{w-1}$  are the number of  $w$ -blocks and  $w - 1$ -blocks. In other words, we define the entropy of the fractional width recoding as the sum of the entropy of the pre-computed table and the entropy of each block. For an  $\mathcal{L}$ -bit secret key recoded with our fractional width method, the average entropy  $\overline{S_F}$  and the worst entropy  $\check{S}_F$  are:

$$\begin{cases} \overline{S}_F = \log_2 |\Omega_B| + \frac{\mathcal{L}}{w-1+p} (pS_w + (1-p)S_{w-1}) \\ \underline{S}_F = \log_2 |\Omega_B| + \min\left(\frac{\mathcal{L}}{w}S_w, \frac{\mathcal{L}}{w-1}S_{w-1}\right) \end{cases} \tag{4}$$

*Entropy of the pre-computed table.* We first evaluate the contribution of the randomized pre-computed table to the entropy of the fractional width method. There are  $k - 2^{w-1}$  exponents randomly chosen in the set of  $2^{w-1}$  elements  $\{2^{w-1} + 1, 2^{w-1} + 2, \dots, 2^w\}$ . Therefore, the entropy of the pre-computed table is:

$$\log_2 |\Omega_B| = \sum_{i=k-2^{w-1}}^{2^{w-1}} \log_2 i - \sum_{i=1}^{2^w-k} \log_2 i. \tag{5}$$

*Entropy of a w-block.* Then, we study the entropy of digits recoded with the width  $w$ . We consider the digits from the upper half table  $x > 2^{w-1}$  first. If  $x > 2^{w-1}$  (probability  $1/2$ ), since  $w$  has been chosen,  $x \in B$ : there are only  $k - 2^{w-1}$  digits in the upper half table. Thus, when  $w$  is selected, the probability of a digit from the upper half table is  $\frac{1}{2(k-2^{w-1})}$ . Next, we consider digits from the lower half table ( $x \leq 2^{w-1}$ ). When  $x \leq 2^{w-1}$  (probability  $1/2$ ), since the lower half of the table is full, there are  $2^{w-1}$  possible digits. Therefore, in the lower half table, the digits have a probability of  $\frac{1}{2 \cdot 2^{w-1}}$ . Clearly, digits in the upper half table have a greater probability than digits in the lower half table. This difference is the origin of entropy loss when  $w$  is selected. More precisely, the entropy of a block of width  $w$  is:

$$S_w = \frac{w+1}{2} + \frac{1}{2} \log_2 (k - 2^{w-1}) \tag{6}$$

Note that in the case of a perfect countermeasure, we have  $S_w = w$  bits.

*Entropy of a w - 1-block.* The case of the width  $w - 1$  is slightly different. When  $x \leq 2^{w-1}$  (probability  $1/2$ ), the choice among the  $2^{w-1}$  digits is uniform. But when  $x > 2^{w-1}$  (probability  $1/2$ ), among the  $2^{w-1}$  possible digits, only  $2^w - k$  can be selected (because  $x \notin B$ ). Therefore,  $k - 2^{w-1}$  digits have a probability of appearance of  $\frac{1}{2} \frac{1}{2^{w-1}}$ , whereas  $2^w - k$  digits have a probability of appearance of  $\frac{1}{2} \left(\frac{1}{2^{w-1}} + \frac{1}{2^w-k}\right)$ . As a consequence, the entropy of a block of width  $w - 1$  is:

$$S_{w-1} = \frac{k-2^{w-1}}{2^w}w - \left(\frac{2^w-k}{2^w} + \frac{1}{2}\right) \log_2 \left(\frac{1}{2^w} + \frac{1}{2^{w+1}-2k}\right) \tag{7}$$

Again, a perfect countermeasure would have  $S_{w-1} = w - 1$  bits.

### 4.3 Consequences on Security

*Simple Power Analysis.* Despite the entropy loss of the unsigned fractional window method, we claim that the representation is still a good SPA countermeasure. Indeed, one can interpret entropy as the equivalent bitlength of the secret when the SPA information can be fully utilized. Therefore, in practice, an SPA countermeasure is at least as strong as its SPA information entropy. As a consequence, if the worst-case entropy is large enough, we know for sure that the

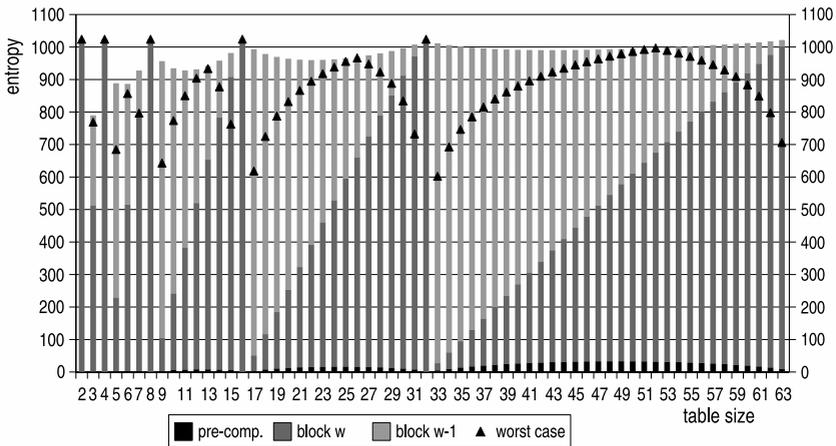


Fig. 1. Average and worst entropy, 1024-bit fractional width method

secret will be safe in any case. More precisely, if the table length  $k$  is carefully chosen, the entropy granted by the unsigned fractional width method is sufficient to thwart all known attacks against partially-compromised RSA secret keys. In particular, we remark that in Fig. 1, for from  $k = 37$  to  $k = 61$ , the worst-case entropy is greater than 800 bits. Additionally, there is currently no attack on RSA or prime field cryptosystems which can take advantage of the bias of the fractional width technique: typically, to be effective, such attacks require the knowledge of the upper or lower bits, or consecutive bits. Thus, we argue that in practice, the security level of the unsigned fractional width technique is even higher than its SPA entropy.

Unfortunately, this is not the case for other exponentiation methods. The binary exponentiation has an entropy of 0 bits, because distinguishing squares from multiplications allows to fully revealing the secret key. In the case of the sliding window method, it is possible to distinguish blocks of consecutive zeros, where no multiplication occurs. Thus, the average entropy of the sliding window representation is  $S_{SW} = \mathcal{L}/2$  bits. And in the worst case, the entropy can go as low as 0 bits. Even if the latter event is unlikely, situations where more than half of the bits are revealed can occur with a relatively high probability.

*Cache Attacks.* Cache analysis allows to distinguish not only different types of operations (such as squares and multiplications), but also reveals information about their operands (in particular, which pre-computed value is accessed by multiplications). On the one hand, standard SPA countermeasures have no effect against the latter type of leakage. On the other hand, thanks to its randomized pre-computed table, our countermeasure makes cache attacks less practical. For instance, with  $k = 53$ , there are more than  $2^{26}$  possible pre-computed tables. And on the top of that, when considering practical cache attacks, one has to take other sources of noise into account: cache misses due to other non-cryptographic

processes, multiple pre-computed values stored in the same cache line, same pre-computed value stored in consecutive cache lines among others.

## 5 Conclusion

We presented a new countermeasure for protecting modular exponentiations against side-channel attacks. Our countermeasure is inspired by the signed fractional width technique, but has an unsigned digit set, which is necessary for achieving high efficiency with RSA or DSA. Because the countermeasure produces a secure recoding of secret values, and because one secret value admits many recodings, our countermeasure protects exponentiations against several types of side-channel attacks: timing attacks, simple power analysis and cache attacks for instance. Our countermeasure is not only highly secure, but is also efficient; unlike other countermeasures which provoke a significant performance drop, thanks to its higher flexibility, our method is about as fast as commonly used exponentiation techniques such as the sliding window method. In fact, in optimal settings and in a practical situation, our method is only 5% slower than the sliding window exponentiation for 512-bit and 1024-bit moduli.

## Acknowledgement

This work was partly supported by National Institute of Information and Communications Technology (NICT).

## References

1. Boneh, D., DeMillo, R., Lipton, R.: On the Importance of Checking Cryptographic Protocols for Faults. In W. Fumy (Ed.): *Advances in Cryptology–Eurocrypt’97*. LNCS **1233**, pp.37-51. Springer-Verlag (1997).
2. Chevallier-Mames, B., Ciet, M., Joye, M.: Low-Cost Solutions for Preventing Simple Side-Channel Analysis: Side-Channel Atomicity. In *IEEE Transactions on Computers*, **53(6)**, pp.760-768. IEEE Computer Society (2004).
3. Coppersmith, D.: Finding a Small Root of a Bivariate Integer Equation; Factoring with High Bits Known. In U. Maurer (Ed.): *Advances in Cryptology–Eurocrypt’96*. LNCS **1070**, pp.178-189. Springer-Verlag (1996).
4. Coron, J.-S.: Resistance Against Differential Power Analysis for Elliptic Curve Cryptosystems. In Ç. Koç and C. Paar (Eds.): *Cryptographic Hardware and Embedded Systems–CHES’99*. LNCS **1717**, pp.292-302. Springer-Verlag (1999).
5. Kocher, P., Jaffe, J., Jun, B.: Differential Power Analysis. In M. Wiener (Ed.): *Advances in Cryptology–Crypto’99*. LNCS **1666**, pp.388-397. Springer-Verlag (1999).
6. Kocher, P.: Timing Attacks on Implementations of Diffie-Hellman, RSA, DSS, and Other Systems. In N. Kobitz (Ed.): *Advances in Cryptology–Crypto’96*. LNCS **1109**, pp.104-113. Springer-Verlag (1996).
7. Möller, B.: Securing Elliptic Curve Point Multiplication against Side-Channel Attacks. In G. Davida, Y. Frankel (Eds.): *Information Security–ISC’01*. LNCS **2200**, pp.324-334. Springer-Verlag (2001).

8. Möller, B.: Improved Techniques for Fast Exponentiation. In P.J. Lee, C.H. Lim (Eds.): Information Security and Cryptology–ICISC’02. LNCS **2587**, pp.298–312. Springer-Verlag (2003).
9. Okeya, K., Takagi, T.: The Width-w NAF Method Provides Small Memory and Fast Elliptic Scalar Multiplications Secure against Side Channel Attacks. In M. Joye (Ed.): Topics in Cryptology–CT-RSA’03. LNCS **2612**, pp.328–342. Springer-Verlag (2003).
10. Okeya, K., Takagi, T.: A More Flexible Countermeasure against Side Channel Attacks Using Window Method. In C. Walter, Ç. Koç, C. Paar (Eds.): Cryptographic Hardware and Embedded Systems–CHES’03. LNCS **2779**, pp.397–410. Springer-Verlag (2003).
11. Okeya, K., Takagi, T., Vuillaume, C.: On the Exact Flexibility of the Flexible Countermeasure Against Side Channel Attacks. In H. Wang, J. Pieprzyk, V. Varadharajan (Eds.): Information Security and Privacy–ACISP’04. LNCS **3108**, pp.466–477. Springer-Verlag (2004).
12. OpenSSL: The Open Source Toolkit for SSL/TLS. [www.openssl.org](http://www.openssl.org)
13. Percival, C.: Cache Missing for Fun and Profit. Technical report, available at [www.daemonology.net/papers/htt.pdf](http://www.daemonology.net/papers/htt.pdf)
14. Shoup, V.: NTL: A Library for Doing Number Theory. [www.shoup.net/ntl/](http://www.shoup.net/ntl/)
15. Walter, C.: MIST, an Efficient, Randomized Exponentiation Algorithm for Resisting Power Analysis. In B. Preenel (Ed.): Topics in Cryptology – CT-RSA’02. LNCS **2271**, pp.53–66. Springer-Verlag (2002).
16. Yen, S.-M., Chen, C.-N., Moon, S.-J. , Ha, J.: Improvement on Ha-Moon Randomized Exponentiation Algorithm. In C. Park, S. Chee (Eds.): Information Security and Cryptology–ICISC’04. LNCS **3506**, pp.154–167. Springer-Verlag (2005).
17. Yen, S.-M., Joye, M.: Checking Before Output May Not Be Enough Against Fault-Based Cryptanalysis. In IEEE Transactions on Computers, **49(9)**, pp.967–970. IEEE Computer Society (2000).