

# Exploiting the Central Reduction in Lattice-Based Cryptography

TOLUN TOSUN<sup>1</sup>, AMIR MORADI<sup>2</sup>, and ERKAY SAVAS<sup>1</sup>

<sup>1</sup>Sabancı University, İstanbul, Türkiye (e-mail: tolungtosun@sabanciuniv.edu, erkays@sabanciuniv.edu)

<sup>2</sup>Technische Universität Darmstadt, Darmstadt, Germany (e-mail: amir.moradi@tu-darmstadt.de)

The work described in this paper has been supported in part by European Union through the Twinning Project 101079319 (acronym enCRYPTON). Views and opinions expressed are those of the authors only and do not necessarily reflect those of the European Union or the European Research Council. Neither the European Union nor the granting authority can be held responsible for them.

**ABSTRACT** This paper questions the side-channel security of **central reduction** technique, which is widely adapted in efficient implementations of Lattice-Based Cryptography (LBC). We show that the central reduction leads to a vulnerability by creating a strong dependency between the power consumption and the sign of sensitive intermediate values. We exploit this dependency by introducing the novel **absolute value** prediction function, which can be employed in higher-order non-profiled multi-query Side-Channel Analysis (SCA) attacks. Our results reveal that – compared to classical reduction algorithms – employing the central reduction scheme leads to a two-orders-of-magnitude decrease in the number of required SCA measurements to exploit secrets of masked implementations. We particularly show that our approach is valid for the prime moduli employed by Kyber and Dilithium, the lattice-based post-quantum algorithms selected by NIST. We practically evaluate our introduced approach by performing second-order non-profiled attacks against an open-source masked implementation of Kyber on an ARM Cortex-M4 micro-processor. In our experiments, we revealed the full secret key of the aforementioned masked implementation with only **250** power traces without any forms of profiling or choosing the ciphertexts.

**INDEX TERMS** Side-Channel Analysis, Correlation Power Analysis, Post-Quantum Cryptography, Kyber, Dilithium, Plantard, Montgomery, Arithmetic Masking, Centered Reduction

## I. INTRODUCTION

Shor’s algorithm [1] violates the security of traditional public-key cryptography including RSA and ECC through quantum computing. As the development of a larger quantum computer in the number of qubits is being reported each year, the quantum threat gradually becomes a reality. On the other hand, NIST’s post-quantum cryptography contest is in the fourth round, with already selected algorithms. Among the winners, the lattice-based algorithms form the majority: Kyber [2], Dilithium [3] and Falcon [4]. Although the post-quantum algorithms can resist quantum computing attacks, special attention should be paid to SCA attacks [5], [6] when these algorithms are implemented in both hardware and software.

The core operation in LBC is the polynomial multiplication. For an efficient implementation, the Number Theoretic Transform (NTT) stands out as an excellent approach. NTT is indeed a special form of the Fast Fourier Transform (FFT) that operates on a discrete space. An important building block that significantly impacts the efficiency of the NTT algorithm is the modular reduction of integers, concerning the arithmetic for coefficients of polynomials. Classical techniques such as the Montgomery reduction [7] and Barrett reduction [8] are already applied to LBC by the existing literature [9]–[13].

The same holds for the relatively new Plantard reduction scheme [14], [15]. One important distinction regarding the integer reduction in LBC compared to the RSA and ECC is the bit-length of the number to be reduced. LBC requires a reduction of relatively smaller numbers, that usually fit into a single computer word. For instance, Kyber employs a 12-bit coefficient modulus and Dilithium a 23-bit one. Moreover, the signed representation of integers over a modulus instead of the classical unsigned representation is more desired in LBC [9]–[12], [15]. That is to make a central reduction (a.k.a *centered* reduction) to the range  $[-q/2, q/2]$  instead of  $[0, q]$  for an odd modulus  $q$ . The Plantard and Montgomery algorithms enable 2-cycle implementation of central reduction on the ARM Cortex-M4 [10], [11], [15] while Plantard is superior since the output of Montgomery reduction requires an additional subtraction or addition to be correct.

### A. MAIN MOTIVATION.

Overall, the central reduction improves the efficiency of LBC implementations; however, it also creates new possibilities for SCA attacks. A well-known fact regarding the processors realized by CMOS technology is that the power consumption has a relatively strong dependency on the Hamming weight

(HW) of the processed data. Respectively, the sign of a number in  $[-q/2, q/2]$  becomes the dominant factor influencing its power consumption considering 2's complement to represent negative numbers. Motivated by this, we explore the characteristics of the central reduction in terms of SCA leakage and exploit them, particularly in the presence of a masking countermeasure.

### B. RELATED WORK ON MASKING LBC.

It is well known that masking countermeasures provide a promising way of mitigating EM-/power-based SCA attacks. Accordingly, the existing literature on masking LBC is already quite rich. Some examples applied on Kyber [16]–[19], Dilithium [20]–[22], Saber (another promising post-quantum lattice-based Key Encapsulation Mechanism (KEM)) [23], [24] and generic lattice-based encryption [25]–[27]. Indeed, masking the polynomial arithmetic is considered trivial and is achieved by simply repeating the operations. For instance, a polynomial multiplication  $a \cdot b$  can be performed through the random shares  $a^0$  and  $a^1$  individually while satisfy  $a = a^0 + a^1$ , instead of accessing  $a$  in plain. On the other hand, masking the non-linear components of algorithms involves specialized techniques. One simple example is the compression operation in Kyber, which aims to identify the interval each coefficient of a given polynomial  $a$  resides. It is easy to see that such an operation cannot be performed independently on the arithmetic shares, as it is the case with polynomial multiplication. A masked implementation of such non-linear operations typically utilizes arithmetic-to-Boolean mask conversion [16]–[18], [23]–[26] while a Boolean masking scheme ensures  $a = a^0 \oplus a^1$ . Overall, the existing work regarding masking LBC aims to achieve provable first- or higher-order security by introducing more efficient solutions for masking the non-trivial parts of the algorithms. Particularly, current masked implementations on the ARM Cortex-M4 [16], [17], [21], [23], [28] inherit the polynomial arithmetic from the well-known pqm4 library [29], known for providing state-of-the-art yet unprotected implementations of post-quantum algorithms. To the best of our knowledge, there has been no work questioning the difficulty of possible SCA attacks on the linear parts as long as the implementations are proven to be secure in the desired security order.

### C. RELATED WORK ON SCA ATTACKS AGAINST LBC.

On the other hand, there exist many SCA attacks in the literature that target implementations of LBC. They can be seen in two distinct classes: profiled SCA attacks such as [30]–[38] and non-profiled SCA attacks such as [39]–[42]. The non-profiled attacks commonly target polynomial multiplication for which a secret polynomial is multiplied with a publicly known polynomial that changes based on the input given to the victim implementation. In [39], the authors show that the performance of non-profiled attacks against the polynomial multiplication directly depends on the employed multiplication algorithm as well as on the parameters such as the coefficient modulus. While for some instances of LBC,

the non-profiled attacks can take a relatively long time to retrieve secret polynomials, acceleration is possible in certain scenarios by collecting more measurements, *e.g.* as shown in [40], [41] for distinct implementations. The authors of [41] particularly focused on so-called *incomplete* NTT, a special case of NTT that is employed in efficient implementations of both Kyber and Dilithium on the ARM Cortex-M4 [29]. Besides, it is shown in [42] that the polynomial multiplication can be also effectively targeted in the case of hardware implementations. Except [41], the aforementioned non-profiled attacks target unprotected implementations, *i.e.* not masked. On the other hand, profiled attacks demand for a stronger adversary model. Particularly, a freely accessible device for profiling that is identical to the victim must be available. Nevertheless, the majority of published SCA attacks on LBC benefit from a profiling phase. To this end, several operations of LBC have been selected as the target of these attacks. Among them, [30]–[32] focus on the NTT transformation. The attack presented in [30] is distinctive by revealing *single-trace* vulnerabilities of LBC. It should be emphasized that if the masks during the profiling phase are known (*i.e.* they are also considered in the profiles), single-trace profiling attacks cannot be avoided by masking, since SCA leakage of a single execution is measured, where the masks can also be revealed via the profiles. The authors of [33] specifically focus on the multiplication of polynomials with small coefficients. The works [34]–[36] present attacks on encoding/decoding functions that transform binary input into a polynomial or vice versa. We should highlight that [35] presents only a message recovery attack that aims to retrieve the decapsulated message rather than the secret key. The target of [37] is the sampling of challenge polynomials, which have a limited number of non-zero coefficients that are also small in magnitude. Table 1 formalizes the above discussion and positions our study among the attacks existing in the literature.

### D. OUR CONTRIBUTIONS.

Below is a list of the contributions we have made in this work.

- To the best of our knowledge, we present the first study in the literature that is particularly developed to effectively exploit the leakage of SCA-protected implementations of LBC without the need for profiling.
- We show that the central reduction techniques that are widely adapted in LBC lead to a source of effectively exploitable SCA leakage. Particularly, information about the sign of arithmetic shares would ease exploiting the leakage and conducting successful key-recovery attacks.
- We show that the employed coefficient modulus as well as the reduction algorithm and the machine word size affect the SCA leakage of masked implementations of LBC, particularly making non-profiled attacks easier to be conducted. For the aforementioned scenarios, we particularly present the so-called optimal correlation and the number of traces required for different noise levels. Our study reveals that SCA attacks on masked implementations with central reduction are significantly more noise-tolerant.

TABLE 1: Qualitative summary of state-of-the-art SCA attacks on implementations of LBC.

Work	Class	Algorithm	Implementation	Masked	Target Function
this work	Non-Profiled	Dilithium <sup>✱</sup> , Kyber	Cortex-M4 <sup>✱</sup>	✓✓	poly. mult.
[39]	Non-Profiled	Kyber, Saber, NTRU <sup>1</sup>	Cortex-M4 <sup>✱</sup>	✗	poly. mult.
[40]	Non-Profiled	Dilithium	Ref. C	✗	poly. mult.
[41]	Non-Profiled	Dilithium, Kyber	Cortex-M4 <sup>✱</sup>	✓	poly. mult.
[42]	Non-Profiled	Dilithium	Hardware	✗	poly. mult.
[30]	Profiled	→	Cortex-M4 <sup>◇</sup>	✓	NTT
[31]	Profiled	Dilithium	Ref. C	✗	NTT
[31]	Profiled	Dilithium	Ref. C <sup>*</sup>	✓	sparse poly. mult.
[32]	Profiled	Kyber	Ref. C	✗	NTT
[32] <sup>⊠</sup>	Profiled	Kyber	Cortex-M4 <sup>✱</sup>	✗	bin. to poly. small
[33]	Profiled	Dilithium	♣	✗	poly.mult
[34]	Profiled	Kyber, Saber	Cortex-M4 <sup>✱</sup>	✓	poly. to bin.
[35] <sup>⊠</sup>	Profiled	Kyber	Cortex-M4 <sup>✱</sup>	✓	bin. to poly.
[36]	Profiled	Dilithium	Ref. C	✗	bin. to poly.
[37] <sup>⊠</sup>	Profiled	Dilithium, NTRU <sup>1</sup> NTRU Prime <sup>2</sup>	Ref. C Cortex-M4 <sup>✱</sup>	✓	small poly. sampling

✓✓ presents a novel technique to tackle masking without the need for profiling

✱ attacks in the simulation

★ from [43], ◇ from [25]

♣ from [16], ♣ from [26]

→ generic to NTT applications in LBC

\* attacks through an implementation submitted for another project [44]

♣ not reported

⊠ challenge polynomial/message-recovery attack

<sup>1</sup> a post-quantum lattice-based KEM [45]

<sup>2</sup> a post-quantum lattice-based KEM [46]

- We introduce a novel prediction function for non-profiled SCA attacks, namely the absolute value prediction function that well predicts the leakage caused by the adaption of central reduction in masked implementations of LBC. We show that our introduced prediction function is optimal.
- We practice our approach against a first-order masked implementation of the lattice-based post-quantum KEM Kyber. We experimentally show that only a few hundred traces are required to successfully mount a key-recovery attack.

## II. BACKGROUND

### A. NOTATIONS

The notations we followed in the paper are as follows.

- Vectors are represented by bold lowercase letters such as  $\mathbf{x}$ , while matrices are represented by bold uppercase letters such as  $\mathbf{X}$ . Polynomials are represented by lowercase regular letters such as  $x$ . Vector-to-vector, matrix-to-vector, and scalar-vector multiplications are denoted by  $\cdot$  while element-wise multiplication of vectors or matrices is denoted by  $\star$ . The  $i$ -th element (coefficient) of a vector (polynomial) is denoted by the subscripts, such as  $x_i$  ( $x_i$ ).
- The central reduction to the range  $[-q/2, q/2]$  is explicitly denoted by  $\text{mod}^{\pm}q$ , while  $\text{mod}q$  denotes the regular modular reduction to the range  $[0, q)$ . We also use  $\text{mod}q$  when the output range is not important such as in high-level representation of algorithms, *i.e.* pseudocodes. The set of unsigned integers in  $[0, q)$  is denoted by  $\mathbb{Z}_q$ .

Accordingly,  $\pm\mathbb{Z}_q$  represents the signed representation of integers modulo  $q$ , namely the integers in the range  $[-q/2, q/2]$ . We assume an odd  $q$  unless the opposite is explicitly stated.

- Random variables are denoted by uppercase letters such as  $X$ .  $\mathcal{P}(\cdot)$  denotes the probability function,  $E[\cdot]$  the expected value function, and  $\mathcal{N}(\mu, \sigma)$  the noise following a Gaussian distribution with mean  $\mu$  and standard deviation  $\sigma$ .
- Unless otherwise stated, the logarithm is base 2, and  $\oplus$  denotes the exclusive OR.
- Shares of variables are represented by superscripts, such as  $X = X^0 + X^1$ .
- $\mathcal{B}(X)$  denotes the number of bits needed to represent the unsigned integer  $X$ .  $\beta$  denotes the machine word size. In this paper, either  $\beta = 16$  or  $\beta = 32$ .
- $\mathcal{W}_{\beta}(X)$  represents the Hamming weight (HW) of a signed integer  $X$  in  $\beta$ -bit 2's complement representation.
- $\mathcal{S}(X) : \mathbb{Z} \rightarrow \{0, 1\}$  returns the non-negativeness (sign) of the integer  $X$ :

$$\mathcal{S}(X) = \begin{cases} 1, & \text{if } X \geq 0 \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

### B. Lattice-Based Cryptography (LBC)

Here, we briefly review lattice-based post-quantum algorithms from an SCA attack perspective. We focus on Kyber [2] and Dilithium [3], which are among the algorithms selected by NIST at the end of the third round of the post-quantum cryptography standardization process. Afterwards,

we discuss the details of the NTT, which is a crucial primitive for efficiently implementing polynomial arithmetic in LBC.

Most of the lattice-based cryptosystems, including Kyber and Dilithium, operate over the ring of polynomials  $\mathcal{R}_q = \mathbb{Z}_q[x]/(X^n + 1)$  that contains polynomials up to degree  $n - 1$  while the coefficients are in  $\mathbb{Z}_q$ . Arithmetic operations in  $\mathcal{R}_q$  are the main building block for implementing LBC, and it is the main target of SCA attacks as well.

### 1) Kyber

Kyber [2] is a lattice-based post-quantum KEM, a variant of the LPR encryption scheme [47]. For all security levels, Kyber employs  $q = 3329$  and  $n = 256$  to instantiate  $\mathcal{R}_q$ . The key pair is generated by the MLWE [48] equation  $\mathbf{t} = \mathbf{A} \cdot \mathbf{s} + \mathbf{e}$ . The vector of polynomials  $\mathbf{e} \in \mathcal{R}_q^k$  is considered as noise and thrown away after the key generation while  $\mathbf{s} \in \mathcal{R}_q^k$  forms the secret key. On the other hand,  $\mathbf{A} \in \mathcal{R}_q^{k \times k}$  and  $\mathbf{t} \in \mathcal{R}_q^k$  are public. The polynomials in both  $\mathbf{s}$  and  $\mathbf{e}$  are *short*, whose coefficients are sampled from the central binomial distribution  $\mathcal{B}_\eta$  with error distribution  $\eta$ . The sensitive operation in a KEM in terms of non-profiled SCA attacks is the decapsulation function where the secret key is involved [39], [41]. The decapsulation in an LPR scheme such as Kyber is quite simple:  $v - \mathbf{s}^T \cdot \mathbf{u}$ , where  $\mathbf{u} \in \mathcal{R}_q^k$  and  $v \in \mathcal{R}$  together form the ciphertext. Related parameters  $k$  and  $\eta$  are chosen depending on the NIST security level as  $\{2, 3, 4\}$  and  $\{2, 4, 2\}$ , respectively.

### 2) Dilithium

Dilithium is a lattice-based post-quantum signature following the Fiat-Shamir scheme with aborts approach [49]. It employs  $q = 2^{23} - 2^{13} + 1 = 8380417$  and  $n = 256$ . The secret-public key pair for Dilithium is generated through the MLWE equation similar to Kyber, *i.e.*  $\mathbf{t} = \mathbf{A} \cdot \mathbf{s}_1 + \mathbf{s}_2$ . Distinctively, both  $\mathbf{s}_1$  and  $\mathbf{s}_2$  are saved as the secret key while the pseudo-randomly generated matrix  $\mathbf{A} \in \mathcal{R}_q^{k \times l}$  and vector  $\mathbf{t} \in \mathcal{R}_q^l$  are public as in Kyber. The coefficients of the secret polynomials in  $\mathbf{s}_1$  and  $\mathbf{s}_2$  are short as well. Specifically, the secret coefficients are sampled uniformly at random in  $[-\eta, \eta]$ . A natural target for a non-profiled SCA attack on Dilithium is the signature function as it involves the secret key [40], [41], [50]. More precisely, the multiplications  $c \cdot \mathbf{s}_1$  and  $c \cdot \mathbf{s}_2$  are targeted, where  $-$  among the outputs of the signature  $-$  the challenge polynomial  $c \in \mathcal{R}_q$  is public and depends on the input message. The parameters  $(k, l)$  are chosen as  $\{(4, 4), (6, 5), (8, 7)\}$  with respect to the security level.  $\eta$  is chosen the same as in Kyber.

### C. Number Theoretic Transform (NTT)

NTT allows efficient multiplication of polynomials in  $\mathcal{R}_q$ . Given two polynomials  $a \in \mathcal{R}_q$  and  $b \in \mathcal{R}_q$ , the NTT multiplication is performed as follows.

$$\text{NTT}^{-1}(\text{NTT}(a) \star \text{NTT}(b)) \quad (2)$$

To simplify the notation, we denote the NTT transformation for polynomials using ‘hat’ for the rest of the paper, *i.e.*

$\hat{a} = \text{NTT}(a)$ . The element-wise multiplication in the NTT domain,  $\hat{a} \star \hat{b}$ , is known as the base multiplication.

NTT is considered as an application of the Chinese Remainder Theorem (CRT) to  $\mathcal{R}_q$ . In case  $q \equiv 1 \pmod{2n}$ , a primitive  $2n$ -th root of unity  $\zeta_{2n} \in \mathbb{Z}_q$  exists for which  $\zeta_{2n}^n \equiv -1 \pmod{q}$ . This setting allows a *complete* NTT over  $\mathcal{R}_q$ , where  $x^n + 1$  can be factored down to the linear factors as  $\prod_{i=0}^{n-1} (x - \zeta_{2n}^{2i+1})$ . The NTT transformation indeed computes the remainder from the division of its input polynomial by  $(x - \zeta_{2n}^{2i+1})$  for each  $i$ , resulting in a vector of  $n$  elements, *i.e.* in  $\mathbb{Z}_q^n$ . Consequently, the base multiplication is performed coefficient wise, *i.e.* a modular multiplication for each  $i$ .

The *forward* and *backward* NTT transformations can be efficiently implemented in  $\log(n)$  steps. Each step is called an NTT *layer*. Indeed, the polynomial is recursively split using so-called *butterfly units* until a linear degree is reached. The forward transformation is usually implemented with *Cooley-Tuckey* (CT) butterflies [51] while the backward transformation is commonly realized using *Gentleman-Sande* (GS) butterflies [52] although it is not a must. For an input pair of coefficients  $a_0$  and  $a_1$ , the CT butterfly computes the output pair by

$$\hat{a}_0 = a_0 - a_1 \cdot \delta, \quad \hat{a}_1 = a_0 + a_1 \cdot \delta, \quad (3)$$

where  $\delta$  is called the *twiddle factor*, a power of  $\zeta_{2n}$ .

### 1) Incomplete NTT

Sometimes, due to performance optimizations or restrictions of the operated ring of polynomials, NTT is not computed for all  $\log(n)$  layers [9], [11], [53]–[55]. This is referred to as *incomplete* NTT. In case the NTT is computed for  $m < \log(n)$  layers,  $\hat{a}$  for  $a \in \mathcal{R}_q$  is a vector with  $2^m$  elements and each element is a polynomial with a degree of  $\log(n) - m$ . Then, the base multiplication refers to the multiplication of degree- $(\log(n) - m)$  polynomials. For instance, Kyber employs  $q = 3329$  and  $n = 256$  allowing a 7-layer NTT while  $\log(n) = 8$ . Therefore, the base multiplication in this setting is achieved by performing 128 individual multiplications of degree-1 polynomials, as demonstrated in Algorithm 1.

---

**Algorithm 1** Base Multiplication for  $(\log(n) - 1)$ -layer incomplete NTT

---

**Input:**  $\hat{s}, \hat{c}$ ; the resulting vectors from  $(\log(n) - 1)$ -layer forward NTT transformation on  $s, c \in \mathcal{R}_q$

**Output:**  $\hat{z} = \hat{s} \star \hat{c}$

- 1: **for**  $\forall i \in \{0, \dots, n/2 - 1\}$  **do** ▷ Compute  $\hat{z}_i = \hat{s}_i \cdot \hat{c}_i$
  - 2:  $\hat{z}_{i,0} \leftarrow \hat{s}_{i,0} \cdot \hat{c}_{i,0} + \hat{s}_{i,1} \cdot \hat{c}_{i,1} \cdot \delta_i \pmod{q}$  ▷  $\delta_i$  is a power of  $\zeta_{2n}$
  - 3:  $\hat{z}_{i,1} \leftarrow \hat{s}_{i,1} \cdot \hat{c}_{i,0} + \hat{s}_{i,0} \cdot \hat{c}_{i,1} \pmod{q}$
  - 4: **end for**
- 

### D. MODULAR ARITHMETIC

As follows, we briefly review the modular reduction techniques that are adapted in LBC. Compared to ECC or RSA,

the modular arithmetic in LBC deals with relatively shorter integers. Additionally, operating with signed integers in modular arithmetic proves to be more efficient in LBC [9]–[12], [15]. The main reason for this is due to the fact that it simply eliminates the need for an extra addition for preventing negativeness in the butterfly units (see Equation (3)). Table 2 summarizes state-of-the-art reduction schemes implemented on the ARM Cortex-M4. It can be seen that the smallest latency in terms of the number of clock cycles is achieved by central reduction techniques [9], [11], [15], whose output range is centered around 0. As a result, in addition to the NTT transformation, central reduction is also preferred to speed up the base multiplication.

### 1) Barrett Reduction

The Barrett reduction was originally proposed in [8]. Its main idea is to subtract a factor of the modulus  $q$  from the number to reduce by approximating the division of the number by  $q$  through a pre-computed factor and shifting. A signed version of Barrett reduction adapted for LBC is proposed in [13]. The input range of the signed Barrett reduction is  $[-\beta/2, \beta/2)$ , and the output range is  $[0, q]$ . A 9-cycle implementation of the signed Barrett reduction for packed integers dedicated to ARM Cortex-M4 is presented in [11]. Packing of integers refers to storing two  $\beta = 16$ -bit integers in a  $2\beta = 32$ -bit register, which is then passed to the packed reduction function. Later, a 6-cycle implementation of Barrett reduction for packed integers was reported in [9], which performs a central reduction with an output range  $[-q/2, q/2]$ .

### 2) Montgomery Reduction

The Montgomery reduction is first proposed in [7]. Similar to the Barrett reduction, it enables a constant time reduction by eliminating the need for division. A signed version of Montgomery reduction is presented by [13], with an input range  $[-q \cdot \beta/2, q \cdot \beta/2)$  and output range  $(-q, q)$ . While a 3-cycle implementation on ARM Cortex-M4 was initially given by [12] for  $\beta = 16$ , the state-of-the-art implementation of Montgomery reduction [10], [11] takes 2 cycles for both  $\beta = 16$  and  $\beta = 32$ . Also, an 8-cycle implementation of Montgomery reduction for packed integers was presented in [11]. We would like to note that a final correction may be required after the signed Montgomery reduction to find the residue in the signed range  $[-q/2, q/2]$  which is the ultimate goal. However, we computationally found out that the correction step is not needed most of the time. In particular for  $q = 3329$  and  $\beta = 16$ , only %0.1 of the corresponding input range requires a correction while the output is already in  $[-q/2, q/2]$  for the rest. Therefore, we study the output range  $[-q/2, q/2]$  for the sake of simplicity for the rest of this paper.

### 3) Plantard Reduction

The Plantard reduction [14] is a more recent algorithm compared to its counterparts, Montgomery and Barrett. While the original Plantard reduction operates on unsigned integers, the authors of [15] proposed an improved version, which operates

on signed integers to be employed in LBC. The output range of the signed version is  $[-q/2, q/2]$ , the same as the state-of-art Barrett reduction. One advantage of the Plantard reduction is that it enables 2-cycle modular multiplication by a constant, outperforming the 3-cycle Montgomery multiplication. The multiplication by a constant is beneficial for implementing the butterfly units during the NTT transformations (see Equation (3)). On the other hand, the improved Plantard reduction also takes 2 cycles on ARM Cortex-M4, the same as Montgomery. However, Plantard’s 2-cycle implementation enables a larger input range and a smaller output range that is desirable. Specifically, the Plantard reduction generates the output in the exact range  $[-q/2, q/2]$ . In other words, it does not require any final correction. This is a significant improvement over the 2-cycle Montgomery reduction whose output range is  $(-q, q)$ . As a side note, packed reduction takes 5 cycles.

## III. NON-PROFILED SCA ATTACK ON NTT MULTIPLICATION

In this section, we present the general outline of a non-profiled power SCA attack on an implementation of a polynomial multiplication in the NTT domain.

### A. ATTACK OUTLINE

#### 1) Leakage Model

Let us first define the assumed leakage function of the target device based on the random variable  $X$  defined in the space  $\mathcal{X}$ , constant scaling factor  $\alpha$ , and noise sampled from a Gaussian distribution with mean  $\mu$  and standard deviation  $\sigma$ , which is independent of  $X$  as

$$\mathcal{L}(X) = \alpha \cdot \mathcal{W}_\beta(X) + \mathcal{N}(\mu, \sigma). \quad (4)$$

$\mathcal{L}(X)$  is commonly used to simulate SCA leakage of microprocessors in the presence of noise when  $X$  is processed.

#### 2) Adversary Model

Consider the base multiplication  $\hat{s} \star \hat{c}$  where  $s \in \mathcal{R}_q$  is a secret polynomial, and  $c \in \mathcal{R}_q$  is a public polynomial. The goal of an adversary is to reveal  $\hat{s}$  through a non-profiled SCA attack where the attacker has access to the power consumption pattern of the underlying device during the computation of  $\hat{s} \star \hat{c}$ . In a non-profiled attack, the attacker samples the leakage of  $\hat{s} \star \hat{c}$  for  $\nu$  distinct computations where  $\hat{c}$  changes for each measurement. The set of samples recorded for each measurement is referred to as a *trace*.

#### 3) Attack on Complete NTT

As previously mentioned, the base multiplication in the NTT domain is performed element wise. Therefore, each  $\hat{s}_i$  can be attacked independently using the knowledge of  $\hat{c}_i$ . To retrieve  $\hat{s}_i$ , a set of hypotheses is made. Each hypothesis is tested by evaluating the *target function*,  $g(\hat{s}_i, \hat{c}_i)$ , thereby statistically comparing the observed leakages (traces). We use the above-defined random variable  $X$  as the sensitive output of the

TABLE 2: Summary of state-of-the-art reduction implementations on ARM Cortex-M4.

Scheme	$\beta$	Input Range	Output Range	Packed	Cycles
Montgomery [11]	16	$[-q \cdot \beta/2, q \cdot \beta/2)$	$(-q, q)$	✗	2
Montgomery [10]	32	$[-q \cdot \beta/2, q \cdot \beta/2)$	$(-q, q)$	✗	2
Montgomery [11]	16	$[-q \cdot \beta/2, q \cdot \beta/2)$	$(-q, q)$	✓	8
Barrett [13]*	16	$[-\beta/2, \beta/2)$	$[0, q]$	✗	3
Barrett [9]	16	$[-\beta/2, \beta/2)$	$[-q/2, q/2]$	✓	6
Plantard [15]	16	$[-q^2 2^{2\alpha'}, q^2 2^{2\alpha'}]^\star$	$[-q/2, q/2]$	✗	2
Plantard [15]	16	$[-q^2 2^{2\alpha'}, q^2 2^{2\alpha'}]^\star$	$[-q/2, q/2]$	✓	5

\* gives the definition of the algorithm but does not present an implementation on ARM Cortex-M4.

✗  $\alpha'$  is a parameter of Plantard reduction that satisfies  $q < 2^{\beta - \alpha' - 1}$ .

target function,  $X = g(\hat{s}_i, \hat{c}_i)$ , as we consider  $\hat{c}_i$  as a random variable changing over measurements. In this study, the *target function* to reveal  $\hat{s}_i$ , is the multiplication  $g(\hat{s}_i, \hat{c}_i) = \hat{s}_i \cdot \hat{c}_i \bmod q$  (or  $\bmod^{\pm} q$  if the reduction is central). Therefore,  $\mathcal{X} = \mathbb{Z}_q$  (or  $\mathcal{X} = \pm \mathbb{Z}_q$ ), and there are  $q$  hypotheses to test. Accordingly, we use  $X' = g(\hat{s}'_i, \hat{c}_i)$  to denote the evaluation of the chosen target function for the hypothesis  $\hat{s}'_i$ . In a Correlation Power Analysis (CPA) attack [6], the output of the *target function* is transformed into *hypothetical leakages* using a *prediction function* with HW function,  $\mathcal{W}_\beta(X')$ , being the most frequently used prediction function. CPA is based on estimating the correlation, *e.g.* by Pearson correlation coefficient, between measured and hypothetical leakages,  $\hat{\rho}(\mathcal{W}_\beta(X'), \mathcal{L}(X))$ . With a sufficient number of traces, the attacker expects to obtain  $\hat{s}_i = \operatorname{argmax}_{\hat{s}'_i} (\hat{\rho}(\mathcal{W}_\beta(X'), \mathcal{L}(X)))$ . Since the attacker may not know which point in time corresponds to  $\mathcal{L}(X)$  in the traces, the procedure is usually repeated for a subset of all samples points.

#### 4) Attack on Incomplete NTT

For  $(\log(n) - 1)$ -level incomplete NTT such as in Kyber,  $q^2$  hypotheses should be examined since  $\hat{s}_i \cdot \hat{c}_i$  is a multiplication between degree-1 polynomials (see Algorithm 1) [39], [41]. In this case, the attacker can use the lower-degree coefficient  $g(\hat{s}_i, \hat{c}_i) = \hat{z}_{i,0}$  for  $\hat{z}_i = \hat{s}_i \cdot \hat{c}_i$ , the higher-degree coefficient  $g(\hat{s}_i, \hat{c}_i) = \hat{z}_{i,1}$ , or both coefficients  $g(\hat{s}_i, \hat{c}_i) = \{\hat{z}_{i,0}, \hat{z}_{i,1}\}$  as the target. If lower- or higher-degree coefficients are used alone,  $\mathcal{X} = \mathbb{Z}_q$ , similar to when NTT is complete. On the other hand,  $\mathcal{X} = \mathbb{Z}_q^2$  if both coefficients are involved in the target function. Since a number and its additive inverse in 2's complement form<sup>1</sup> are correlated, the number of hypotheses can be reduced to  $q/2$  and  $q^2/2$  for complete and incomplete NTT, respectively. In this manner, the attacker learns either the actual secret or its additive inverse.

#### 5) Application to Kyber and Dilithium

It is important to emphasize that the adversary model as well as the target function explained above can directly be applied to both Kyber and Dilithium. More precisely, performing an

attack on the leakage of  $\hat{s} \star \hat{c}$  to reveal  $\hat{s}$  corresponds to targeting  $s^T \cdot u$  in the case of Kyber and targeting  $c \cdot s_1$  or  $c \cdot s_2$  for Dilithium (see Section II-B).

## B. SECOND-ORDER ATTACK

### 1) Masking

The most promising way to defeat the above-explained SCA attack is masking [16]–[27]. Indeed, masking of the polynomial multiplication is straightforward from an algorithmic perspective, as it can be seen as a linear operation. For a uniformly and randomly generated share  $s^0$ , one computes the other share as  $s^1 = s - s^0$ . Then, the computation  $\hat{s} \star \hat{c}$  is performed on the shares as  $\hat{s}^0 \star \hat{c}$  and  $\hat{s}^1 \star \hat{c}$ . Notice that,  $\hat{s} \star \hat{c} = \hat{s}^0 \star \hat{c} + \hat{s}^1 \star \hat{c}$ . Accordingly,  $\hat{s}_i \cdot \hat{c}_i = \hat{s}_i^0 \cdot \hat{c}_i + \hat{s}_i^1 \cdot \hat{c}_i$  for all  $i$ . This is referred to as *arithmetic masking*. In particular, the order of masking is defined by the number shares representing the secrets. Here in the given example, first-order masking is applied as two shares are used. Since  $\hat{s}_i \cdot \hat{c}_i$  is not computed in plain, the leakage of every single point in SCA traces is expected to be independent of  $\hat{s}_i$  and hence independent of the estimated hypothetical leakage, *e.g.*  $\mathcal{W}_\beta(X')$ . It is noteworthy to mention that an assumed condition for such a claim is that each share  $s^0$  and  $s^1$  should individually follow a uniform distribution.

### 2) Combination Function

In order to conduct successful CPA attacks on a first-order masked implementation, the leakage associated to two shares should be combined using a pre-processing function. Since the shares are processed individually (not simultaneously), their associated leakages appear in different time samples. Intuitively for the application studied in this work, the attacker combines the observed leakages associated to  $X^0 = g(\hat{s}_i^0, \hat{c}_i)$  and  $X^1 = g(\hat{s}_i^1, \hat{c}_i)$  while mean-free product is known as the most efficient pre-processing (combination) function [56], [57]. For leakages of the random shares,  $\mathcal{L}(X^0)$  and  $\mathcal{L}(X^1)$ , the mean-free product is defined as

$$\mathcal{C}(\mathcal{L}(X^0), \mathcal{L}(X^1)) = (\mathcal{L}(X^0) - E[\mathcal{L}(X^0)]) \cdot (\mathcal{L}(X^1) - E[\mathcal{L}(X^1)]) \quad (5)$$

The mean-free product is the most preferred combination function because it effectively reduces noise by removing

<sup>1</sup>Negative integers in 2's complement form are represented by inverting all bits of the corresponding positive integer and adding 1 to the result.

constant offsets, thereby enhancing the signal-to-noise ratio and emphasizing relevant variations. Needless to say,  $E[\mathcal{L}(X^0)]$  and  $E[\mathcal{L}(X^1)]$  are approximated by the sample means over the trace set. For sake of completeness, we include the absolute difference combination function [58] in Appendix , which is also an efficient and frequently used combination function.

### 3) Optimal Prediction Function (OPF)

[56] shows that an optimal prediction function must be computed in order to efficiently perform a second-order CPA attack by means of a combination function

$$f_{opt}(x) = E[\mathcal{C}(\mathcal{L}(X^0), \mathcal{L}(X^1)) | X = x]. \quad (6)$$

Recall that  $X = X^0 + X^1 \pmod q$  (or  $\pmod{\pm q}$ ). In simple words, the expected value of the combination function is calculated for the given unmasked value which can be hypothetically computed. A CPA attack performed on a masked implementation by means of such a pre-processing function is referred to as *Higher-Order CPA (HOCPA)*. Similar to the first-order attack, the attacker estimates the correlation  $\hat{\rho}(f_{opt}(X'), \mathcal{C}(\mathcal{L}(X^0), \mathcal{L}(X^1)))$  to rank the hypotheses. The correlation achieved by  $f_{opt}$  and the correct hypothesis is referred to as the *optimal correlation*, denoted by  $\hat{\rho}_{opt} = \hat{\rho}(f_{opt}(X), \mathcal{C}(\mathcal{L}(X^0), \mathcal{L}(X^1)))$ . In this study, we use HOCPA with the mean-free product as the distinguisher and discuss the optimal prediction function and its efficiency for the presented adversary model in different reduction scenarios. However, we replicate some of the experiments discussed in the subsequent sections for the absolute difference combination function.

### 4) Conditioning on Zero-Values

As the chosen target function is a multiplication for this study, one can fine-tune the optimal prediction function by considering zero-value public data, *i.e.*  $\hat{c}_i = 0$ . Notice that for the complete NTT,  $X = 0$  if and only if  $\hat{c}_i = 0$  (assuming  $\hat{s}_i \neq 0$ ), because of the speciality of 0 in multiplication. Then,  $X^0$  and  $X^1$  become 0 as  $X^j = \hat{s}_i^j \cdot 0$  and  $f_{opt}(0) = E[\mathcal{C}(Y_0, Y_1) | X^0 = 0, X^1 = 0] = \alpha^2 \cdot (E[\mathcal{W}_\beta(X) + \mathcal{N}(\mu, \sigma)])^2$ , see Equation (4). When the NTT is incomplete,  $X$  can be 0 even though  $\hat{c}_i$  is non-zero, because the target  $X$  is the addition of two multiplications which can sum up to 0. Recall that  $X$  is a function of  $\hat{z}_{i,j}$  in this case, which is formulated in Algorithm 1. However, since having  $\hat{c}_i = 0$  is less likely for the incomplete NTT ( $1/q^2$  for  $\hat{c}_{i,0} = \hat{c}_{i,1} = 0$ , assuming uniformly random  $\hat{c}_i$ ) we do not concentrate on this case.

## IV. LEAKAGE OF SIGNED INTEGERS MODULO $q$

In this section, we study the distribution of HW of the signed representation of integers modulo  $q$ , *i.e.* the effect of central reduction on HW. Accordingly, we compare signed and unsigned arithmetic in terms of SCA leakage. We evaluate the primes that are employed in Kyber and Dilithium with  $q = 3329$  and  $q = 8380417$ , respectively, and study the carrier primes that are employed for Dilithium to perform short

polynomial arithmetic [9], namely  $q = 257$  for Dilithium2 and Dilithium5, and  $q = 769$  for Dilithium3. We should note that masking the short polynomials in their range is possible [59]. One important factor for computing the HW of negative integers is the machine word size  $\beta$  which does not have any effect on the HW of positive integers. The machine word size is usually  $\beta = 16$  for  $q = 257$ ,  $q = 769$ , and  $q = 3329$  while  $\beta = 32$  for  $q = 8380417$  in software implementations. Unless otherwise stated, we take these values for  $\beta$  in the studied adversary model. However, we discuss the role of distinct values of  $\beta$  in the SCA leakages.

### A. HW AS A SIGN INDICATOR

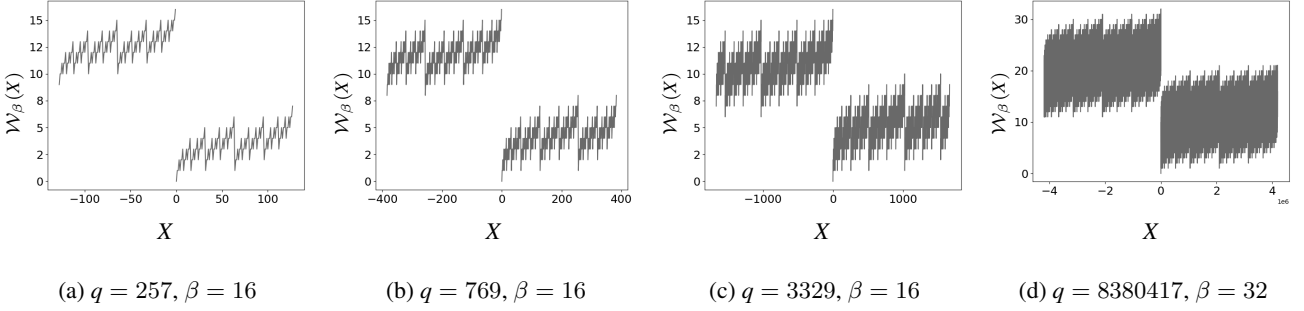
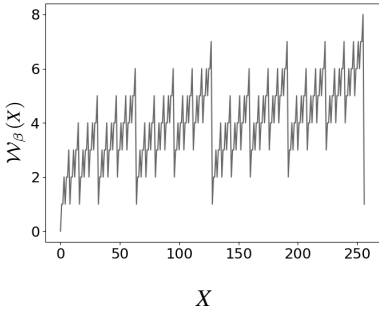
The main observation that led to this study is the clear separation of HW of the non-negative side of  $\pm\mathbb{Z}_q$ , namely  $[0, q/2]$  and the negative side  $[-q/2, 0)$ . As an intuition, consider  $q = 257$ ; the positive interval of  $\pm\mathbb{Z}_{257}$  corresponds to  $[0, 128]$ , for which the maximum HW is  $\mathcal{W}_\beta(127) = 7$ . In other words, the HW of integers  $[0, 128]$  lies in  $[0, 7]$ . On the other hand, the negative side of  $\pm\mathbb{Z}_{257}$  corresponds to  $[-128, 0)$ , where the HWs are in the range of  $[9, 16]$  assuming 2's complement representation with machine word size  $\beta = 16$ . Consequently, the HW of a number in  $\pm\mathbb{Z}_{257}$ , reveals its sign immediately. Figure 1 visualizes our observation for all the primes analyzed in this study. Note that there is an overlap between the HW ranges  $[-q/2, 0)$  and  $[0, q/2]$ , for  $q = 769$ ,  $q = 3329$ , and  $q = 8380417$ . For instance, the HW of non-negative integers in  $\pm\mathbb{Z}_{3329}$  are distributed in  $[0, 10]$  while that of negative integers are in the range  $[6, 16]$ . Therefore, there is an overlap for five possible HWs in the interval  $[6, 10]$  out of a total of 17 possible values assuming the machine word size  $\beta = 16$ .

Based on our observation in Figure 1, we write the following equality for  $\mathcal{W}_\beta(X)$  where  $X \in \pm\mathbb{Z}_q$ .

$$\mathcal{W}_\beta(X) = \mathcal{S}(X) \cdot \gamma + (1 - \mathcal{S}(X)) \cdot (\beta - \gamma) + e, \quad (7)$$

for some inner-cluster error term  $e \in \mathcal{E}$  where  $E[e] = 0$  and  $0 < \gamma < \beta/2$ . Note that  $\gamma$  stands for the mean of HWs given  $X$  is non-negative, *i.e.*  $\gamma = E[\mathcal{W}_\beta(X) | \mathcal{S}(X) = 1]$ . Similarly, let  $\gamma^-$  denote the mean HW given  $X$  is negative,  $\gamma^- = E[\mathcal{W}_\beta(X) | \mathcal{S}(X) = 0]$ . Table 3 demonstrates the values of  $\gamma$  and  $\gamma^-$  depending on  $q$ . Note that  $\gamma^- \approx \beta - \gamma$ , which allows us to simplify Equation (7). Naturally,  $\gamma^-$  approaches  $\gamma$  as the so-called gap  $\mathcal{M}(q, \beta) = \beta - \log(q)$  decreases. On the other hand, the expected value  $E[|e|]$  is not affected by  $\beta$  while it slightly increases as  $q$  gets larger. Additionally, whether the integers modulo  $q$  are represented in signed or unsigned form has no impact on this expected value.

We would like to note that the argument made in this section does not apply to the unsigned representation of integers modulo  $q$ , namely  $\mathbb{Z}_q$ . It can be seen in Figure 2 that no clear ranges can be identified for unsigned integers when observing their HW. In LBC, the upper half of  $\mathbb{Z}_q$  is considered as negative, namely  $[-q/2 + q, q - 1]$ . However, Table 3 shows that the mean of HW of negative side, denoted by  $\gamma^*$ , is approximately  $\gamma + 1$ , independent of  $\beta$  and  $q$ .

FIGURE 1: Distribution of HW of integers in  $[-q/2, q/2]$  in 2's complement representation.FIGURE 2: Distribution of HW of integers in  $[0, q)$  for  $q = 257$ .

$q$	257	769	3329	8380417
$\gamma$	3.48	4.16	5.19	10.99
$\gamma^-$	$\beta - 3.5$	$\beta - 4.17$	$\beta - 5.2$	$\beta - 11$
$\gamma^*$	4.5	5.16	6.19	11.99
$E[ e ]$	1.33	1.4	1.55	2.34

TABLE 3: The mean and standard deviation of HW for positive and negative integers with distinct  $q$ .  $\gamma$  denotes the mean of HWs for positive integers in  $^{\pm}\mathbb{Z}_q$ , namely  $[0, q/2]$ .  $\gamma^-$  denotes the mean of HWs for negative integers in  $^{\pm}\mathbb{Z}_q$ , namely  $[-q/2, -1]$ .  $\gamma^*$  denotes the mean of HWs for  $[q/2 + 1, q - 1]$ , the set of integers in  $\mathbb{Z}_q$  which are considered to be negative in LBC, and  $e$  is the error term defined in Equation (7).

## B. IMPACT OF SIGNED ARITHMETIC ON OPTIMAL CORRELATION

To formally assess the impact of signed arithmetic on SCA leakages, we compare the optimal correlation achieved by the state-of-the-art combination function, mean-free product, between the cases when the modular reduction is central and when it is non-central. For consistency and comparison purposes, we also provide correlation results for the unprotected scenario. Figure 3 presents the estimated (optimal) correlation for distinct prime  $q$  and machine word size  $\beta$  and reduction scenarios. It can be seen that  $\hat{\rho}_{opt}$  estimated for central reduction achieves more than twice of the one estimated for the non-central reduction, particularly for  $\beta = 16$  and  $\beta = 32$ . Indeed,  $\hat{\rho}_{opt}$  is an increasing function of  $\beta$  for a given  $q$  when the reduction is central, complying with our

initial observation in Table 3. As  $\beta$  decreases,  $\hat{\rho}_{opt}$  for the central reduction reaches that of the non-central case as the limit. More importantly, the correlation shows a strong resistance to noise for the signed case. For instance, when  $\sigma = 5$  and  $q = 257$  and  $\beta = 16$  (such as a masked software implementation of Dilithium),  $\hat{\rho}_{opt}$  reaches  $\hat{\rho} = 0.24$  while we observe  $\hat{\rho} = 0.017$  for the unsigned case. We should refer to Table 2 showing that  $\beta$  increases the input range of the reduction algorithms. However, our analysis shows that it further increases the associated SCA leakages. It also makes sense to compare the optimal correlation achieved in case of Boolean masking with the other results. Similar to what presented in [56], the correlation for Boolean masking reaches  $\hat{\rho} = 0.35$  for an 8-bit implementation in a noiseless scenario. Similar to the non-central reduction case, this drops rapidly with the noise, e.g. to  $\hat{\rho} = 0.02$  for  $\sigma = 5$ .

The same pattern is observed in the unprotected scenario. Notably, the correlation is significantly greater when the reduction is central beyond certain values of  $\sigma$ . In particular for  $\sigma = 4$  and  $q = 3329$ , the correlation reaches  $\hat{\rho} = 0.62$  with the signed representation whereas it is  $\hat{\rho} = 0.37$  with the unsigned representation. However, it is important to highlight that, for the same level of  $\sigma$  and  $q$ , the central reduction increases the optimal correlation by  $\approx 64\times$ . More generally, as illustrated in Figure 3, the impact of signed arithmetic on SCA leakage is more prominent in the protected case. An interesting result is that the optimal correlation with the central reduction is relatively close to the correlation observed in unprotected case with non-central reduction, depending on the margin  $\mathcal{M}(q, \beta)$ . Specifically, for  $q = 257$  and  $\sigma \geq 3$ , the difference in correlation between the masked unsigned and unprotected signed cases is less than 0.05. This leads to the conclusion that the effect of masking can diminish when central reduction is applied.

We replicated the same analysis for the absolute difference combination function [58], as detailed in Figure 12 (in Appendix). Although the estimated optimal correlations are slightly lower across all values of  $q$ ,  $\beta$ , and  $\sigma$ , the impact of  $\beta$  and  $\sigma$  follows the same pattern observed for the mean-free product.

For this evaluation, we computed  $f_{opt}$  as a look-up table for each case. This has a time complexity of  $O(q^2)$  to create  $f_{opt}$ , which can be costly for large  $q$ , such as  $q = 8380417$ . In



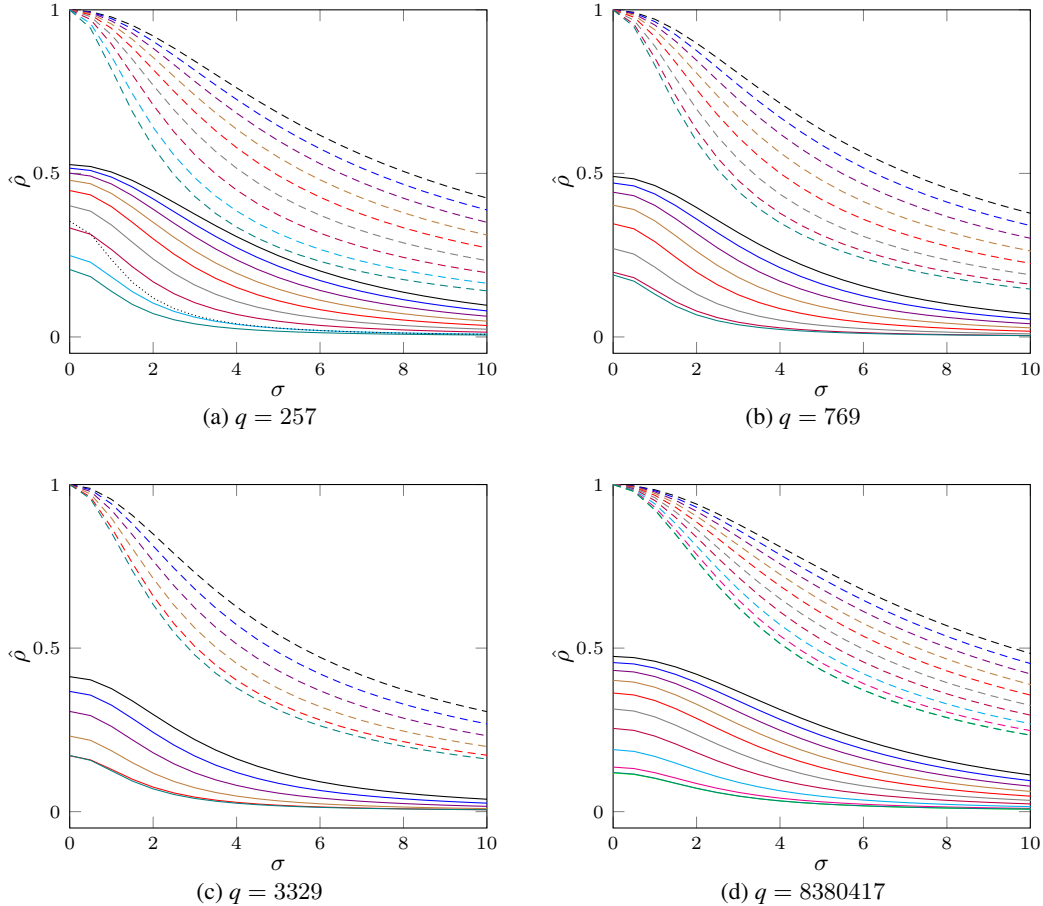


FIGURE 3: Correlation with respect to the noise standard deviation  $\sigma$  for different  $q$  and  $\beta$ , reduction strategies and protection scenarios. For the first-order masked case, estimations are performed with 1 million samples uniformly taken for  $X^0$  and  $X^1$  and optimal correlation  $\hat{\rho}_{opt}$  is reported in y-axis. Otherwise, estimations are performed with the same number of samples uniformly taken for  $X$  and  $\hat{\rho}(\mathcal{W}_\beta(X), \mathcal{L}(X))$  is reported.

- ★ Masked (solid), unprotected (dashed).
- ★ Reduction to  $[-q/2, q/2]$  for  $q = 257$ ,  $q = 769$  and  $q = 3329$ :  $\beta = 16$  (black),  $\beta = 15$  (blue),  $\beta = 14$  (violet),  $\beta = 13$  (brown),  $\beta = 12$  (red),  $\beta = 11$  (gray),  $\beta = 10$  (purple),  $\beta = 9$  (cyan).
- ★ Reduction to  $[-q/2, q/2]$  for  $q = 8380417$ :  $\beta = 32$  (black),  $\beta = 31$  (blue),  $\beta = 30$  (violet),  $\beta = 29$  (brown),  $\beta = 28$  (red),  $\beta = 27$  (gray),  $\beta = 26$  (purple),  $\beta = 25$  (cyan),  $\beta = 24$  (magenta),  $\beta = 23$  (green).
- ★ Reduction to  $[0, q)$  (teal).
- ★ Boolean masking only in (a) for  $\beta = 8$  and  $q = 256$  (dotted).

Section V, we deal with explicit formulas for  $f_{opt}$  dedicated to central reduction.

### C. INFORMATION THEORETIC ANALYSIS

Additionally, we investigate the Mutual Information (MI) between  $X$  and the HW leakage. In particular, we compute  $\mathcal{I}(X, \mathcal{W}_\beta(X))$  for the unprotected case, while in the first-order masked case, we compute the multivariate MI, denoted by  $\mathcal{I}(X, \mathcal{W}_\beta(X^0), \mathcal{W}_\beta(X^1))$ . Figure 4 presents numerical results for the MI, accross different values of  $q$ ,  $\beta$  and reduction

scenarios<sup>2</sup>. Aligning with our observations from the previous section, MI between  $X$  and the HW is higher when the signed representation is employed, which depends on the gap  $\mathcal{M}(q, \beta)$ . This result holds true for both masked and unprotected cases. In the unprotected scenario, the increase in MI is nearly 1, reflecting the leakage introduced by the sign. On the other hand, with masking, the increase in MI is up to  $\approx 4\times$  with  $\beta$ . Recall that the correlation-based analysis also indicates that negative impact of central reduction is more pronounced in the masked case. Note that, MI is a employed

<sup>2</sup>MI results for  $q = 8380417$  with masking enabled are not provided due to the computational complexity of the experiments. Nevertheless, they can be inferred from the results of the other primes studied.

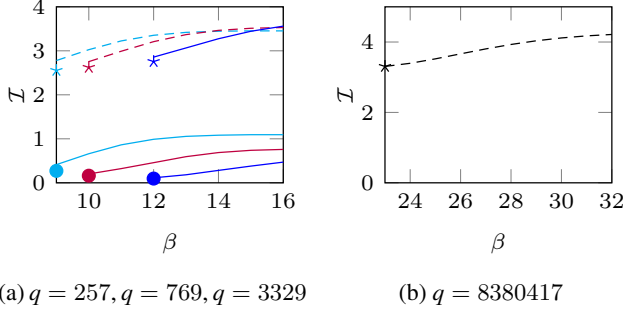


FIGURE 4: Mutual information for different  $q$ ,  $\beta$ , reduction strategies, and protection scenarios. For the unprotected case,  $\mathcal{I}(X, \mathcal{W}_\beta(X))$  is reported while for the first-order masked case  $\mathcal{I}(X, \mathcal{W}_\beta(X^0), \mathcal{W}_\beta(X^1))$  is reported in y-axis.

- ★ Reduction to  $[-q/2, q/2]$ : masked (solid), unprotected (dashed). Reduction to  $[0, q]$ : masked (circle), unprotected (star).
- ★ Moduli:  $q = 257$  (cyan),  $q = 769$  (purple),  $q = 3329$  (blue).  $q = 8380417$  (black).

$\beta \backslash q$	257	769	3329	$\beta \backslash q$	8380417
16	-0.741	-0.732	-0.661	32	-0.706
15	-0.723	-0.706	-0.607	31	-0.684
14	-0.698	-0.671	-0.527	30	-0.655
13	-0.665	-0.621	-0.405	29	-0.618
12	-0.618	-0.545	-0.223	28	-0.570
11	-0.548	-0.427		27	-0.507
10	-0.442	-0.241		26	-0.424
9	-0.273			25	-0.314
				24	-0.173
				23	-0.003

TABLE 4: Estimations of  $\hat{\rho}(\mathcal{W}_\beta(X), f_{opt}(X))$  for different moduli  $q$  and  $\beta$ . The estimations are performed with 1 million uniformly random samples for  $X$  while the reduction is central.

in side-channel analysis [60], as a generic distinguisher where the attacker is not required to predict the device's leakage model. Given that the leakage model in our case is linear with HW as defined in Equation (4), and results with MI results are consistent with our previous findings, we stick with CPA for the remainder of the study.

## V. ABSOLUTE VALUE PREDICTION FUNCTION

When the target operation is protected using Boolean masking and the underlying circuit is a noisy Hamming weight of intermediates (as in Equation (4)), it is shown by [56] that  $\mathcal{W}_\beta$  can be effectively used as the optimal prediction function for HOCPA attacks. However, this is not necessarily the case when masking is arithmetic. To provide an intuition, we estimated  $\hat{\rho}(\mathcal{W}_\beta(X), f_{opt}(X))$ , as proposed in [56], to measure the accuracy of prediction functions for the studied adversary model, presented in Table 4. The results indicate a significant correlation loss in all scenarios. Hence, in this section,

we search for an explicit formula for the optimal prediction function in case of arithmetic masking. Precisely, we show that the absolute value function can be used as the optimal prediction function when targeting arithmetic masking where central reduction is employed.

### A. DISTRIBUTION OF THE SECRET KNOWING THE SIGN OF SHARES

Consider two uniformly random variables  $X^0, X^1 \in \pm \mathbb{Z}_q$  and their modular addition  $X^0 + X^1 \bmod^\pm q$ . Given the signs of both variables, Figure 5 demonstrates the probability distribution of  $X^0 + X^1 \bmod^\pm q$ . As depicted in the figure, there are two cases for the distribution given the sign of both random variables. If the sign of  $X^0$  and  $X^1$  are the same, namely  $\mathcal{S}(X^0) = \mathcal{S}(X^1)$ , then the probability is distributed around  $\pm q/2$ . Otherwise, it is centered around 0. Indeed, the distributions correspond to the convolution of probability distribution functions. More precisely, one of

$$\mathcal{P}(X^0 = x^0 | X^0 < 0), \quad \mathcal{P}(X^0 = x^0 | X^0 \geq 0)$$

is convoluted to one of

$$\mathcal{P}(X^1 = x^1 | X^1 < 0), \quad \mathcal{P}(X^1 = x^1 | X^1 \geq 0).$$

Now suppose that  $X^0$  and  $X^1$  are arithmetic shares representing a secret intermediate variable  $X = (X^0, X^1)$ . Then, the above discussion shows that information about the sign of the individual shares leads to a strong effect on the distribution of the secret  $X$ .

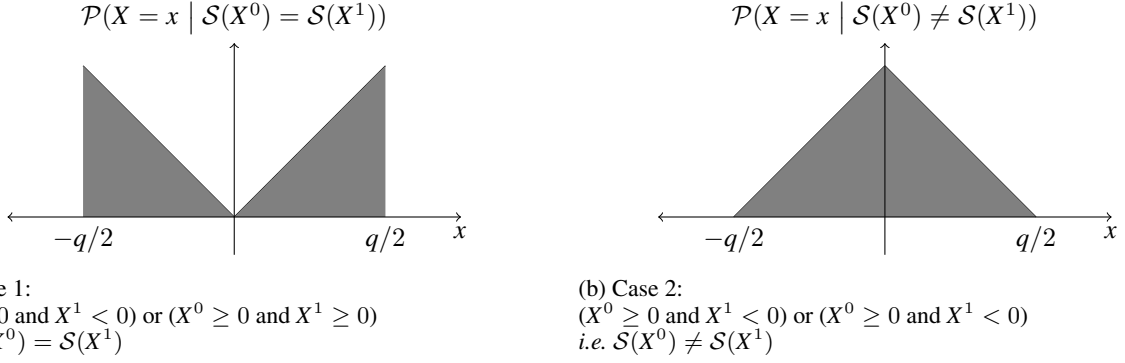
### B. A MODEL FOR MEAN-FREE PRODUCT

In the previous section, we showed that the HW of 2's complement representation of an integer in  $[-q/2, q/2]$  is a noisy indicator of its sign. Also, recall the leakage in CMOS circuits which is highly relevant to the HW of processed data (see Equation (4)). Now, let  $Y_0 = \mathcal{L}(X^0)$  and  $Y_1 = \mathcal{L}(X^1)$  denote the leakage associated to the random shares  $X^0$  and  $X^1$ . The mean-free product can be written as follows.

$$\begin{aligned} \mathcal{C}(Y_0, Y_1) = & \left( \alpha_0 \cdot \mathcal{W}_\beta(X^0) + \mathcal{N}(\mu_0, \sigma_0) - \mathbb{E}[\alpha_0 \cdot \mathcal{W}_\beta(X^0) + \mathcal{N}(\mu_0, \sigma_0)] \right) \cdot \\ & \left( \alpha_1 \cdot \mathcal{W}_\beta(X^1) + \mathcal{N}(\mu_1, \sigma_1) - \mathbb{E}[\alpha_1 \cdot \mathcal{W}_\beta(X^1) + \mathcal{N}(\mu_1, \sigma_1)] \right) \end{aligned} \quad (8)$$

For the sake of simplicity, we assume  $\alpha = \alpha_0 = \alpha_1$ ,  $\mu = \mu_0 = \mu_1$ , and  $\sigma = \sigma_0 = \sigma_1$ . As  $X^0$  and  $X^1$  are uniformly random signed integers in  $\pm \mathbb{Z}_q$  and represented by  $\beta$  bits in the computer memory,  $\mathbb{E}[X^0] = \mathbb{E}[X^1] = \beta/2$ . Then, we can write

$$\begin{aligned} \mathcal{C}(Y_0, Y_1) = & \left( \alpha \cdot \mathcal{W}_\beta(X^0) + \mathcal{N}(\mu, \sigma) - (\alpha \cdot \beta/2 + \mu) \right) \cdot \\ & \left( \alpha \cdot \mathcal{W}_\beta(X^1) + \mathcal{N}(\mu, \sigma) - (\alpha \cdot \beta/2 + \mu) \right) \end{aligned} \quad (9)$$


 FIGURE 5: Probability distributions of  $X = X^0 + X^1 \bmod^{\pm} q$  for  $X^0, X^1 \in \pm \mathbb{Z}_q$ 

and by distributing the terms and as  $\mathcal{N}(0, \sigma) = \alpha \cdot \mathcal{N}(0, \sigma/\alpha)$ ,

$$\begin{aligned} \mathcal{C}(Y_0, Y_1) = \alpha^2 & \left( \mathcal{W}_{\beta}(X^0) - \beta/2 + \mathcal{N}(0, \sigma/\alpha) \right) \cdot \\ & \left( \mathcal{W}_{\beta}(X^1) - \beta/2 + \mathcal{N}(0, \sigma/\alpha) \right) \end{aligned} \quad (10)$$

By plugging Equation (7) into Equation (10) we have

$$\begin{aligned} \mathcal{C}(Y_0, Y_1) = \alpha^2 & \left( \mathcal{S}(X^0) \cdot (\gamma - \beta/2) + (1 - \mathcal{S}(X^0)) \cdot \right. \\ & \left. (\beta/2 - \gamma) + e_0 + \mathcal{N}(0, \sigma/\alpha) \right) \cdot \\ & \left( \mathcal{S}(X^1) \cdot (\gamma - \beta/2) + (1 - \mathcal{S}(X^1)) \cdot \right. \\ & \left. (\beta/2 - \gamma) + e_1 + \mathcal{N}(0, \sigma/\alpha) \right) \end{aligned} \quad (11)$$

### C. CONDITIONAL PROBABILITY OF SIGN EQUALITY

As explained above and shown by Figure 5, we conclude that

$$\begin{aligned} \mathcal{P}(X = x \mid \mathcal{S}(X^0) = \mathcal{S}(X^1)) &= (2/q) \cdot (|x|/(q/2)) \\ &= |x| \cdot (4/q^2). \end{aligned} \quad (12)$$

Based on the dependency of  $X$  on  $\mathcal{S}(X^0)$  and  $\mathcal{S}(X^1)$ , we show that the conditional probability  $\mathcal{P}(\mathcal{S}(X^0) = \mathcal{S}(X^1) \mid X = x)$  is a multiple of  $|x|$ .

$$\begin{aligned} \mathcal{P}(\mathcal{S}(X^0) = \mathcal{S}(X^1) \mid X = x) &= \\ \frac{\mathcal{P}(X = x \mid \mathcal{S}(X^0) = \mathcal{S}(X^1)) \cdot \mathcal{P}(\mathcal{S}(X^0) = \mathcal{S}(X^1))}{\mathcal{P}(X = x)} \\ &= \frac{|x| \cdot (4/q^2) \cdot 1/2}{1/q} = (2/q) \cdot |x| \end{aligned} \quad (13)$$

### D. ESTIMATING THE OPTIMAL PREDICTION FUNCTION

We make use of the conditional probability to formally estimate  $E[\mathcal{C}(Y_0, Y_1) \mid X = x]$  as

$$\begin{aligned} & E[\mathcal{C}(Y_0, Y_1) \mid X = x, \mathcal{S}(X^0) = \mathcal{S}(X^1)] \cdot \\ & \mathcal{P}(\mathcal{S}(X^0) = \mathcal{S}(X^1) \mid X = x) + \\ & E[\mathcal{C}(Y_0, Y_1) \mid X = x, \mathcal{S}(X^0) \neq \mathcal{S}(X^1)] \cdot \\ & \mathcal{P}(\mathcal{S}(X^0) \neq \mathcal{S}(X^1) \mid X = x) \end{aligned} \quad (14)$$

Considering the terms including  $e_0$  and  $e_1$  as error, we can write  $E[\mathcal{C}(Y_0, Y_1) \mid X = x]$  as

$$\begin{aligned} & (\gamma - \beta/2)^2 \cdot (2/q) \cdot |x| + \\ & (\gamma - \beta/2) \cdot (\beta/2 - \gamma) \cdot (1 - 2/q \cdot |x|) + e_c \\ & = (\gamma - \beta/2)^2 \cdot (4/q) \cdot |x| - (\gamma - \beta/2)^2 + e_c \end{aligned} \quad (15)$$

where

$$\begin{aligned} e_c &= E[e_0 \cdot e_1 \mid X = x] + 2E[e_1 \cdot \mathcal{S}(X^0) (\gamma - \beta/2) + \\ & (1 - \mathcal{S}(X^0)) (\beta/2 - \gamma) \mid X = x] \\ &= E[e_0 \cdot e_1 \mid X = x] + \\ & 4(\gamma - \beta/2) \cdot E[e_0 \cdot \mathcal{S}(X^1) \mid X = x] \end{aligned} \quad (16)$$

Note that  $e_c$  can be derived for any valid  $\beta$  when  $E[e_0 \cdot e_1 \mid X = x]$  and  $E[e_0 \cdot \mathcal{S}(X^1) \mid X = x]$  are pre-computed. This approach is beneficial particularly if  $q$  is large (*e.g.* 8380417) and the attacker does not know  $\beta$  in advance. Intuitively,  $e_c$  is relatively small and its impact decreases as the margin  $\mathcal{M}(q, \beta)$  increases. With sufficient  $\mathcal{M}(q, \beta)$ , Equation (15) is accurately approximated by

$$E[\mathcal{C}(Y_0, Y_1) \mid X = x] \approx c_0 \cdot |x| + c_1, \quad (17)$$

where  $c_0 = (\gamma - \beta/2)^2 \cdot 4/q$  and  $c_1 = -(\gamma - \beta/2)^2$ . Figure 6 visualizes these estimations and the corresponding error for two distinct cases of  $q$  and  $\beta$ . Since constants  $c_0$  and  $c_1$  do not affect the correlation,  $f_{abs}(x) = |x|$  can be used as the optimal prediction function.

### E. ON THE ACCURACY OF $F_{ABS}$

Since Equation (16) does not allow to give an exact explicit formula for  $e_c$ , we proposed approximating  $f_{opt}(X)$  using the absolute value function. In order to evaluate the accuracy loss as a result of using  $f_{abs}$  instead of Equation (15), we estimated  $\hat{\rho}(f_{abs}(X), f_{opt}(X))$ . Table 5 presents the estimations, indicating  $\hat{\rho} > 0.99$  for all studied  $q$  and  $\beta$  couples mentioned in Section IV. Therefore, we conclude that the absolute value prediction function is highly accurate for the attacks we consider in this work. However, for any other settings ( $q$

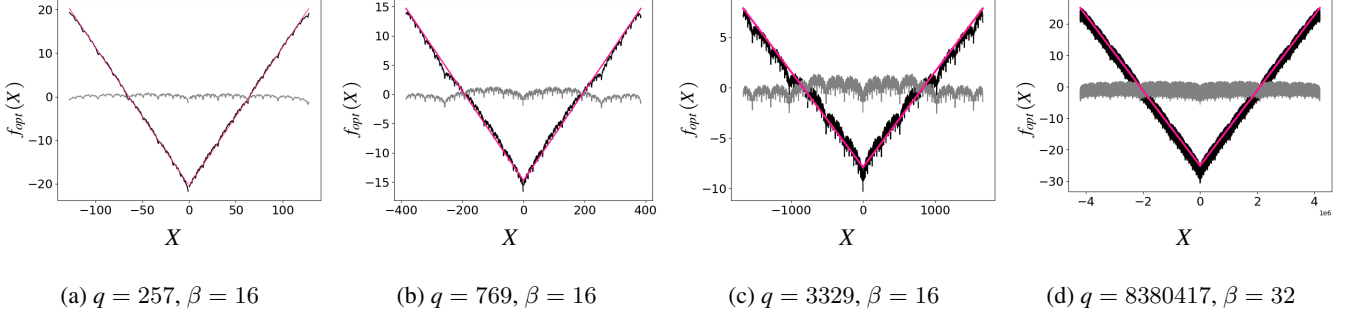


FIGURE 6: Visualization of optimal prediction function for the central reduction range  $[-q/2, q/2]$ .  $f_{opt}(X)$  (black),  $e_c$  (gray),  $c_0 \cdot |X| + c_1$  (pink).

$\beta \backslash q$	257	769	3329	$\beta \backslash q$	8380417
16	<b>0.999</b>	<b>0.997</b>	<b>0.992</b>	32	<b>0.998</b>
15	0.999	0.996	0.984	31	0.998
14	0.998	0.993	0.961	30	0.997
13	0.997	0.987	0.877	29	0.995
12	0.994	0.970	0.546	28	0.992
11	0.985	0.908		27	0.983
10	0.956	0.635		26	0.961
9	0.825			25	0.892
				24	0.659
				23	0.214

TABLE 5: Estimations of  $\hat{\rho}(f_{abs}(X), f_{opt}(X))$  for different moduli  $q$  and  $\beta$ . The estimations are performed with 1 million uniformly random samples for  $X$  while the reduction is central.

and  $\beta$  couples), where the absolute value prediction function leads to an undesired accuracy, the  $f_{opt}(X)$  function can be computed for all possible values of  $X$  as  $E[\mathcal{C}(Y_0, Y_1) | X]$  (e.g. Equation (15)).

The  $f_{abs}$  can also be used with the absolute difference combination function [58]. The corresponding results for this configuration are presented in Appendix .

### F. ALTERNATIVE PREDICTION FUNCTION

To take the advantage of zero-value public data as explained in Section III-B, we define the alternative absolute value prediction function as follows.

$$f_{abs}^*(X) = \begin{cases} ((\beta/2)^2 - c_1)/c_0, & \text{if } X = 0 \\ |X| & \text{otherwise} \end{cases} \quad (18)$$

Recall that  $E[\mathcal{C}(Y_0, Y_1) | X^0 = 0, X^1 = 0] = (\beta/2)^2$ .

## VI. SIMULATION RESULTS

In this section, we present the result of HOCPA attacks explained in Section III making use of simulated traces. In our simulations, we consider different noise levels and various reduction scenarios for each of the studied  $q$ . We particularly compare the required number of traces with and without central reduction. Needless to say that  $f_{abs}$  and  $f_{abs}^*$  are used as the prediction function when the reduction is central while

$f_{opt}$  is computed as a look-up table otherwise (as done in Section IV-B).

### A. SIMULATED TRACES

We perform the following routine to generate simulated traces for base multiplication.

- 1) Generate a secret key vector  $\mathbf{s} \in^{\pm} \mathbb{Z}_q^m$  uniformly at random.
- 2) Generate a public vector  $\mathbf{c} \in^{\pm} \mathbb{Z}_q^m$  uniformly at random.
- 3) If masking is enabled, sample  $\mathbf{s}^0 \in^{\pm} \mathbb{Z}_q^m$  uniformly at random. Then, apply first-order masking to  $\mathbf{s}$  such that  $\mathbf{s} = \mathbf{s}^0 + \mathbf{s}^1 \bmod^{\pm} q$  (or  $\bmod q$  depending on the reduction scheme). Else if masking is disabled, set  $\mathbf{s}_i^0 = 0 \forall i \in \{0, \dots, m-1\}$  and  $\mathbf{s}^1 = \mathbf{s}$ .
- 4) Compute  $\mathbf{s}^0 \star \mathbf{c}$  and  $\mathbf{s}^1 \star \mathbf{c}$ . For each  $i \in \{0, \dots, m-1\}$  and  $j \in \{0, 1\}$ , perform a modular multiplication with (or without) central reduction as  $\mathbf{z}_i^j \equiv \mathbf{s}_i^j \cdot \mathbf{c}_i \bmod^{\pm} q$  (or  $\bmod q$ ) with  $m$  being the number of elements in  $\mathbf{s}$  and  $\mathbf{c}$ . Particularly, the modular reduction is performed in  $[-q/2, q/2]$  or  $[0, q]$  to simulate the leakage for different reduction scheme presented in Table 2.
- 5) Compute the HW of each  $\mathbf{z}_i^j$  as the simulated leakage ( $2m$  individual results).
- 6) Apply Gaussian noise to the simulated leakages as in Equation (4), with  $\alpha = 1, \mu = 0$ , and the given standard deviation  $\sigma$ .
- 7) Go back to Step 2 until  $\nu$  traces are generated.

For computing the HWs, 16-bit or 32-bit 2's complement representations are used depending on  $q$ , as explained in Section IV.  $m, \nu, \sigma, q$  as well as the reduction range (Step 4) are pre-defined parameters. We set  $m = 100$  for all simulations. Note that the traces generated by this routine simulate a base multiplication for a *complete* NTT transformation, where the element-wise multiplication is just a modular multiplication. While  $q = 8380417$  allows a complete NTT with the ring dimension  $n = 256$ , other moduli  $q = 257, q = 769$ , and  $q = 3329$  do not allow complete NTT but allow incomplete NTT of 7 layers. Recall that the element-wise multiplication is the multiplication of degree-1 polynomials in that case. However, the simulations aim to benchmark our introduced absolute value prediction function and compare leakage of

different reduction schemes. Therefore, there is no harm in doing the simulations as if the NTT is complete, which only affects the number of hypotheses. The comparison between the hypothetical and observed leakages is not affected by this behavior. Note also that we use  $s$  and  $c$  here instead of the notation  $\hat{s}$  and  $\hat{c}$  given in Section III.

## B. EVALUATION

Figure 7 presents the corresponding results, with respect to  $q$ ,  $\nu$ , and  $\sigma$ , while the success rate refers to ( $\#$  correctly predicted  $s_i/m$ ). It should be noted that the number of traces is displayed on a logarithmic scale. As evident by the results, the implementations with central reduction are significantly more vulnerable to these non-profiled HOCPA attacks in terms of the number of traces. For instance, when  $q = 3329$  and  $\sigma = 0$ , the attack against non-central reduction needs 1500 traces to succeed, which is  $6\times$  more than the number of traces needed when the reduction is central. Moreover, the noise  $\sigma$  has a greater impact on the attacks on implementations with non-central reduction compared to the central case. When  $q = 3329$  and  $\sigma = 4$ , the attack on non-central reduction needs 45 k traces to succeed, which is  $31\times$  more than what is required in case of central reduction. The difference with respect to the number of traces reaches  $123\times$  for  $q = 257$  and  $\sigma = 4$ . We conclude that HOCPA with  $f_{abs}$  targeting central reduction remains a major threat in different noise levels conforming with the observation shown in Figure 3. Based on the aforementioned decrease in the number of traces required to attack, employing central reduction in masked LBC might be not the best choice from the SCA perspective. Although central reduction is sometimes preferred for efficiency purposes, non-central reduction can harden higher-order SCA attacks in security-demanding applications.

In general, the attack on central reduction needs relatively small number of traces to succeed. However, the number of traces for a successful attack depends on the margin  $\mathcal{M}(q, \beta)$ , which is slightly worse for  $q = 3329$  compared to the other primes considered here. For  $q = 3329$  and  $\sigma = 4$ , HOCPA with  $f_{abs}^*$  requires only 1400 traces to succeed. In a noiseless scenario, where  $q = 257$  and  $\sigma = 0$ , the attack only needs around 110 traces. As previously stated, we consider only two cases  $\beta = 16$  and  $\beta = 32$ . However, we anticipate from Figure 3 that, as  $\beta$  decreases, the number of required traces to attack the central reduction schemes gets closer to that when a non-central reduction scheme is employed. As anticipated, the advantage of using  $f_{abs}^*$  over  $f_{abs}$  highly depends on the number of times when  $c_i = 0$ . For instance when  $q = 257$  and  $\sigma = 4$ ,  $f_{abs}^*$  leads to  $\%13$  reduction in the number of required traces in our experiments. As the chance of observing zero-values decreases when  $q$  increases, the advantage of  $f_{abs}^*$  decreases as well. For example, when  $q = 3329$  and  $\sigma = 4$ ,  $f_{abs}^*$  reduces  $\nu$  by  $\%9$  compared to  $f_{abs}$ .

The disadvantage of signed arithmetic is also evident in the unprotected but noisy scenario, where the number of required traces increases by up to two orders of magnitude.

For example, when  $q = 257$  and  $\sigma = 4$ , the number of traces required with the central reduction is  $9.5\times$  higher than the non-central reduction case. Notice that the increase in  $\nu$  when using signed arithmetic is smaller compared to the protected case, by a factor of  $13\times$  for  $q = 257$ . Consequently, the simulation results align with our observation in Section IV, indicating that the negative impact of signed arithmetic is more significant in masked schemes. Also observe the statistical results for first-order attacks on unsigned arithmetic and second-order attacks on signed arithmetic are closely aligned. In the extreme case of  $q = 257$ , the number of required traces for these two attacks, despite their different orders, is nearly identical. For  $q = 8380417$ , the difference in  $\nu$  for a successful attack between these two scenarios is only  $2.8\times$ .

## VII. PRACTICAL RESULTS: APPLICATION TO KYBER

In this section, we present the result of applying our proposed approach to perform successful HOCPA attacks on a protected implementation of Kyber. It is important to note that we have previously compared the SCA leakage between central and non-central reduction. Therefore, the motivation of this section is to evaluate the difficulty of conducting successful HOCPA attacks targeting central reduction using real data. Source code of the implementations and the attack scripts as well as the simulations presented in the previous section are publicly available.<sup>3</sup>

### A. TARGET IMPLEMENTATION

We focus on the ARM Cortex-M4 specific open-source and first-order masked implementation of Kyber from [26]<sup>4</sup>. The polynomial arithmetic of the implementation is mostly in assembly, ported from the pqm4 project [43] and employs the Montgomery reduction that we illustrated in Section II-D. We also created a second version of the victim implementation by integrating the latest iteration of polynomial arithmetic from pqm4 which employs the Plantard reduction based on [15]<sup>5</sup>. Hereafter, we denote the untouched target implementation with Montgomery reduction by  $\Psi_M$  and the in-house version with Plantard reduction by  $\Psi_P$ . In particular, we focus on the function `basemul_asm` which implements the base multiplication  $\hat{s} * \hat{c}$  in the incomplete NTT domain for both Montgomery and Plantard versions. We should note that – to the best of our knowledge – all masked implementations of post-quantum algorithms on the ARM Cortex-M4 that have been reported in the literature are built on top of pqm4 by directly porting the linear operations including polynomial arithmetic [16], [17], [21], [23], [28]. Therefore, we believe that assessing the most recent iteration of pqm4, featuring state-of-the-art polynomial arithmetic, would be beneficial. The open-source Kyber implementation [26] employs the Montgomery reduction since the more efficient Plantard reduction did not exist when the polynomial implementation

<sup>3</sup><https://github.com/toluntosun21/ExploitingCentralReduction>

<sup>4</sup>[https://github.com/uclcrypto/pqm4\\_masked/](https://github.com/uclcrypto/pqm4_masked/) commit hash: **5fe90ba**

<sup>5</sup><https://github.com/mupq/pqm4> commit hash: **3743a66**

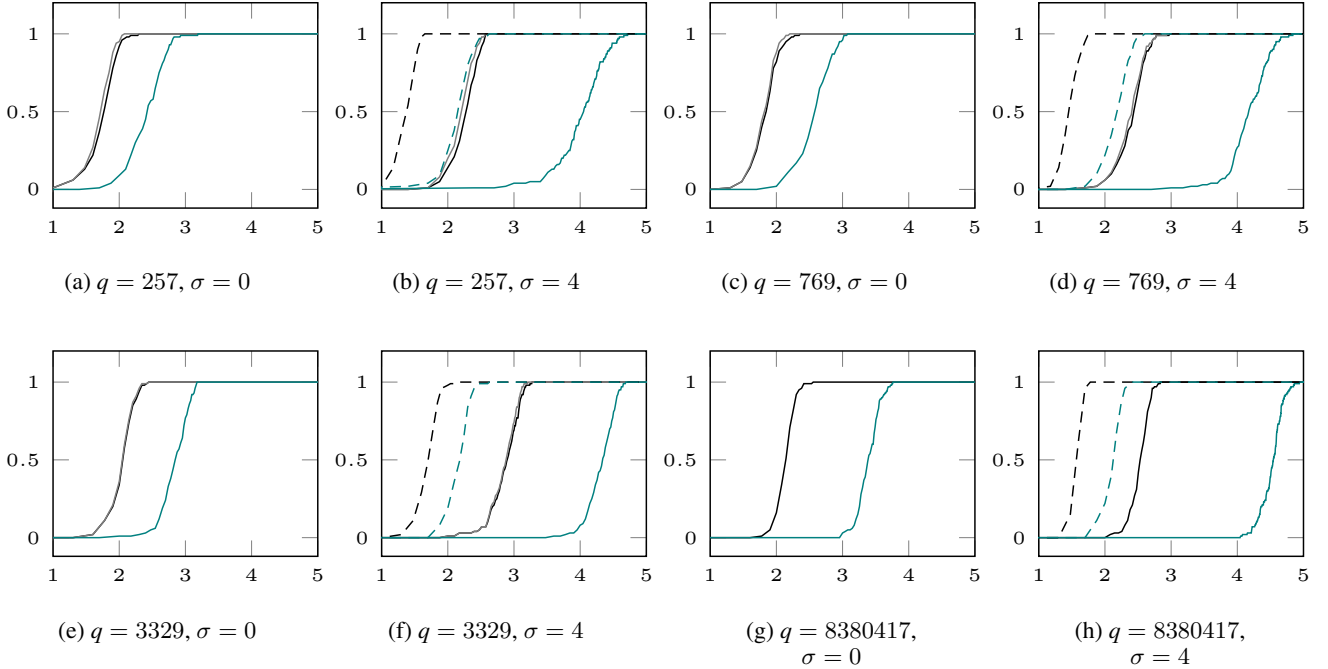


FIGURE 7: Success rates (in y-axis) of first-order and second-order CPA attacks on simulated traces.  $x$ -axis denotes  $\log_{10} \nu$ . For each point in the curves, 100 experiments have been performed with random data.

- ★ SCA Protection: first-order masked (solid), unprotected (dashed).
- ★ Reduction ranges:  $[-q/2, q/2]$  (black, gray),  $[0, q]$  (teal).
- ★ Prediction functions for HOCPA:  $f_{abs}$  (black),  $f_{abs}^*$  (gray),  $f_{opt}$  (teal).

was imported from pqm4. Our experiments are centered around the medium security level, *i.e.* Kyber768, though it does not affect our approach and results.

## B. SETUP

We used NewAE ChipWhisperer CW308 UFO board to collect power traces. The victim program was running on a STM32F303, which is equipped with an ARM Cortex-M4. The frequency of the core is set to 7.3 MHz by an external reference clock which is also given to the power-collecting facility (analog-to-digital converter) while 4 power samples are recorded at each clock cycle. We provided a trigger signal for the power-collecting module to indicate the beginning of the function `basemul_asm` for the first share. Hence, only the samples related to the base multiplication were recorded. The attacks have been performed using the `scared` library<sup>6</sup>, with an in-house developed Python model that mimics the intended Kyber implementation. A laptop equipped with an AMD Ryzen™ 7 7840HS<sup>7</sup> 8-core processor and 64 GB RAM was used for running the attack.

## C. ATTACK DETAILS

A mean trace over 1000 traces is presented in Figure 8. It should be noted that the iterations of the function

`basemul_asm` are visible through the mean trace for both shares. In order to reveal each  $\hat{s}_i$ , in the corresponding attacks we have only taken into account the relevant part of the power traces based on the iterations. We used a constant offset to combine the leakages associated to two shares (by mean-free product as explained in Section III-B) based on the pattern observed in Figure 8. It is noteworthy to mention that the same strategy can be easily adapted via educated guesses without prior knowledge of the specific implementation. Recall that  $\hat{s}_i$  is a degree-1 polynomial in Kyber, and two coefficients must be predicted together based on the outline presented in Section III. We tested  $q \cdot q/2$  hypotheses ( $\approx 2^{22.4}$  as  $q = 3329$ ) with  $f_{abs}$  as the prediction function, so that either the actual secret or its additive inverse is found (for both  $\Psi_M$  and  $\Psi_P$ ). The target of the attack is the higher-degree coefficient of each  $\hat{s}_i \cdot \hat{c}_i$ , precisely  $g(\hat{s}_i \cdot \hat{c}_i) = \hat{z}_{i,1}$  computed in Line 2 of Algorithm 1. When  $f_{abs}$  is used as the prediction function, an hypothesis and its additive inverse,  $\pm \hat{s}_i$ , gets the same correlation score due to the nature of absolute value function.<sup>8</sup>

<sup>8</sup>One option to distinguish the correct hypothesis from its additive inverse is to re-run the same HOCPA attack on two hypotheses  $\pm \hat{s}_i$  using the sign function  $\mathcal{S}$ , see Equation (1). This intuition is based on the fact that not every bit of the intermediate values equally contributes in the amount of power consumption, *i.e.* not an ideal HW model. This leads  $f_{opt}$  to be not fully symmetric with respect to the  $y$  axis, see Figure 6.

<sup>6</sup><https://pypi.org/project/scared/>

<sup>7</sup><https://www.amd.com/en/products/processors/laptop/ryzen/7000-series/amd-ryzen-7-7840hs.html>

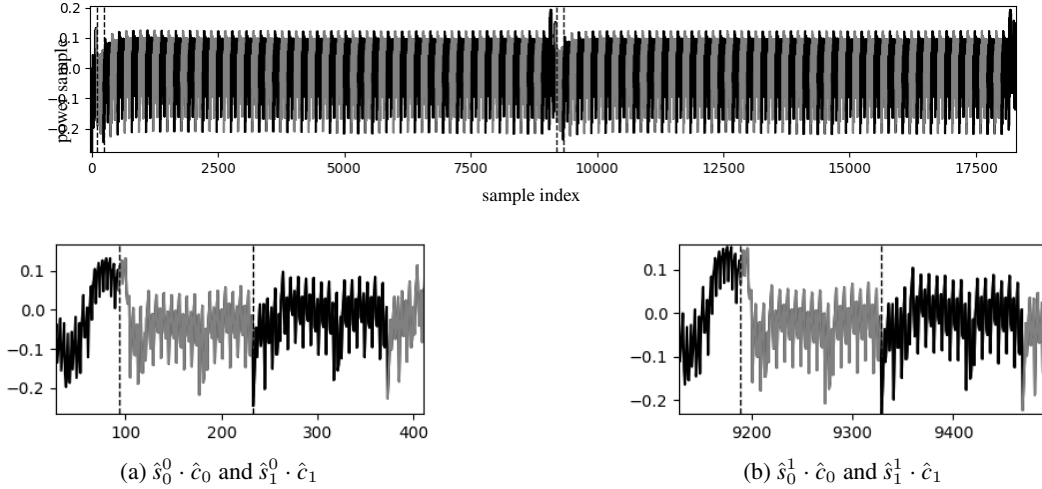


FIGURE 8: The mean power trace associated with the execution of the base multiplication function `basemul_asm` in  $\Psi_M$  for both shares,  $\hat{s}_0 \star \hat{c}$  and  $\hat{s}_1 \star \hat{c}$ . The iterations of the function are marked by interleaving black and gray colors. Due to loop unrolling, 64 iterations are observed for each share instead of  $n/2 = 128$ . Recall that the ring dimension  $n = 256$  for Kyber and 7-layer NTT is performed. The first iterations of `basemul_asm` for both shares are marked and zoomed in (a) and (b). The same observation applies to  $\Psi_P$ .

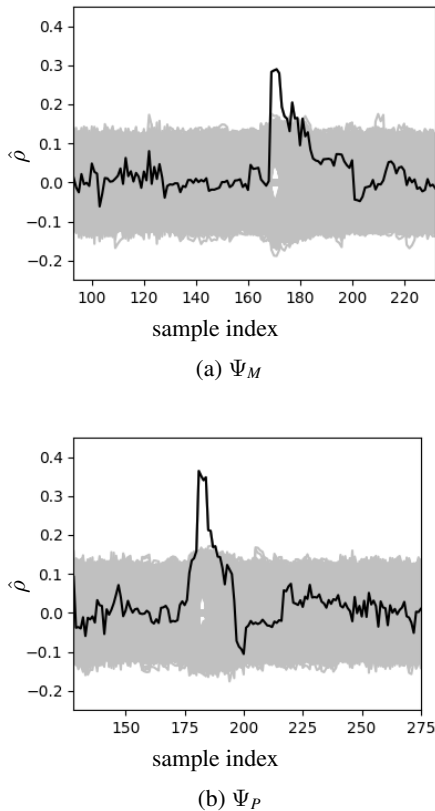


FIGURE 9: HOCPA with  $\nu = 1000$  and  $f_{abs}$  targeting  $\hat{s}_0$  for both  $\Psi_M$  and  $\Psi_P$ . The correlation scores of incorrect hypotheses are in gray, and for the correct hypothesis in black.

#### D. EVALUATIONS

Let us start the evaluations by exemplary presenting the result of the individual attacks on  $\hat{s}_0$  for both  $\Psi_M$  and  $\Psi_P$  in Figure 9. The correlation peaks for the correct hypotheses are observed in the corresponding time samples for the secret coefficient. Observe that the correlation for the correct hypotheses are around 0.3, which can be considered a major correlation for a higher-order attack. Needless to say, the correlation coefficient changes for different values of  $i$ . From a better perspective, Figure 10a and Figure 10b present the efficiency of our introduced prediction function  $f_{abs}$  in terms of the number of traces needed to succeed. Consistent with the simulation results,  $f_{abs}$  is very effective against arithmetic masking with central reduction. In particular, the attacks require 850 and 550 traces to fully recover the secret of the evaluated implementations. The reason why the attack on  $\Psi_M$  requires more traces to succeed is that the secret coefficients  $\hat{s}_i$  for even values of  $i$  lead to lower correlation scores in general compared to the rest of the attack. While this observation can be micro-architecture and implementation specific, we did not concentrate on improving it as the overall attack still leads to a reasonably low number of traces. We should also remark that the aim of this study is not to compare  $\Psi_M$  and  $\Psi_P$  since both implementations employ central reduction; rather the goal is to show that the approach generalizes to central reduction techniques.

As a reference, we also included the success rate of a classical first-order CPA by the  $\mathcal{W}_\beta$  prediction function performed on the same but unprotected implementations in Figure 10c and Figure 10d. In order to keep the consistency, we used the same part of the power traces as those considered in HOCAs. Distinctively, we used both lower- and higher-degree coefficients from the output of each  $\hat{s}_i \cdot \hat{c}_i$  as the target

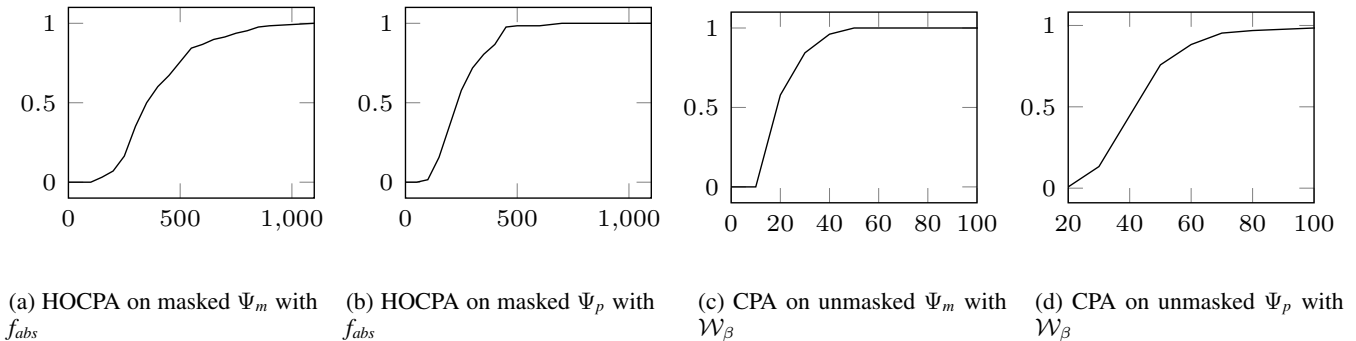


FIGURE 10: Success rates (in y-axis) of the SCA attacks against Kyber. x-axis denotes  $\nu$ . The success rate refers to (# correctly predicted  $\hat{s}_i/128$ ). Retrieving  $\pm\hat{s}_i$  is considered as success.

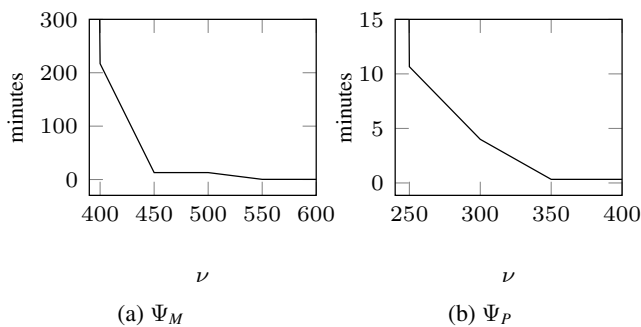


FIGURE 11: Time required for the lattice attacks to successfully retrieve the whole secret polynomial  $s$  using the predictions for  $\hat{s}_i$  which are obtained by HOCPAs.

function, namely  $g(\hat{s}_i, \hat{c}_i) = \{\hat{z}_{i,0}, \hat{z}_{i,1}\}$  (see Algorithm 1). We should also note that  $f_{abs}$  is designed to work with a single coefficient and we leave construction of a prediction function which takes multiple coefficients as the future work.

#### E. BONUS: COMBINING WITH LATTICE ATTACK

As given in [61], the attacker can take the advantage of the fact that SCA attacks are performed in the NTT domain. Intuitively, there exists  $3329^{256}$  possibilities for the NTT domain secret  $\hat{s}$  while this number is  $5^{256}$  for the normal domain secret  $s$ , indicating an over-determined system. The authors of [61] created an LWE instance from the inverse NTT transformation and showed that by retrieving only 38 out of 128 pairs of NTT domain coefficients, the rest of the coefficients can be revealed by solving the LWE, which is practically solvable. The challenge is to assess which NTT domain coefficients was predicted correctly. We applied a series of lattice attacks until an instance is successfully solved, starting from the subset of predictions to the NTT domain coefficients with highest correlation scores. Each execution of the lattice attack takes 20 seconds, returning success in case of the included subset of coefficients were correctly retrieved by the preceding SCA attack. Figure 11 depicts the relationship between the number of traces and the time needed by the

lattice attack to succeed in our experiments. In particular for  $\Psi_M$ , the attack succeeds after 643 trials ( $\approx 3.5$  hours) with 400 traces while it succeeds with 250 traces for  $\Psi_P$  after 27 trials ( $\approx 10$  minutes)<sup>9</sup>. Indeed,  $\nu = 250$  is considered as a very small number to conduct a successful non-profiling higher-order SCA attack on a masked implementation. The lattice attack is implemented using fpylll library<sup>10</sup>.

#### VIII. CONCLUSIONS AND FUTURE WORK

In this paper, we investigated the impact of various reduction schemes on SCA leakage of the implementation of modular arithmetics, with a specific focus on the central reduction. State-of-the-art masking LBC, *e.g.* [16]–[18], [20], [21], [23], [25]–[28], concentrated on developing gadgets to handle non-linear operations and considered the linear part of the algorithms such as the polynomial arithmetic as relatively trivial to mask due to its transparency to arithmetic masking (*i.e.* repeating the operation on each share individually). However, our study reveals that the design decisions such as the reduction technique for the linear parts have a significant impact on the exploitability of the associated SCA leakages and hence on the number of traces required for a successful attack. Our study exposes this fact by showing that a HOCPA attack becomes significantly easier when the masked polynomial multiplication  $\hat{s} \star \hat{c}$  is targeted and the modular reduction in the victim’s implementation is central. Specifically, the use of signed arithmetic notably increases SCA leakage. We quantify this leakage by studying the optimal correlation and the number of traces required for a successful HOCPA attack.

Our findings reveal that the signed representation of integers modulo  $q$  leads to a strong dependency between the sign of an integer and its HW in  $2^s$ ’s complement form. We assessed this dependency through simulations involving the parameter sets employed by the post-quantum cryptography winners Kyber and Dilithium. We also efficiently exploited this source of leakage, by introducing the absolute value prediction function. We believe that our work is unique in the literature as it is the only *non-profiled* SCA attack particularly designed

<sup>9</sup>Except for  $\nu = 250$  for  $\Psi_P$  and  $\nu = 400$  for  $\Psi_M$ , the timing of lattice attacks are approximations.

<sup>10</sup><https://pypi.org/project/fpylll/>



to efficiently exploit higher-order leakages. We further have showcased our approach targeting a first-order masked implementation of Kyber. As our attack does not require profiling and is successful with only 250 traces (in our experiments and using our measurement setup, we claim that utilization of the central reduction in masked implementations indeed may ease SCA attacks. To the best of our knowledge, we report the lowest number of traces for a successful non-profiled second-order SCA attack against masked implementations of LBC. We leave the generalization of our introduced absolute value prediction function to higher orders as a work for the future.

As another outcome of our study, it can also be concluded that finding the sample points in power traces associated with the random arithmetic shares is trivial in masked implementations of LBC. Consequently, it is highly recommended to apply shuffling on top of arithmetic masking. Indeed, it is relatively easy to shuffle the base multiplication since  $\hat{s}_i \cdot \hat{c}_i$  are performed independently for each  $i$  leading to  $\frac{n}{2}!$  possible permutations.

**REFERENCES**

[1] P. W. Shor, “Algorithms for quantum computation: Discrete logarithms and factoring,” in *35th FOCS*. IEEE Computer Society Press, Nov. 1994, pp. 124–134.

[2] P. Schwabe, R. Avanzi, J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, G. Seiler, D. Stehlé, and J. Ding, “CRYSTALS-KYBER,” National Institute of Standards and Technology, Tech. Rep., 2022, available at <https://csrc.nist.gov/Projects/post-quantum-cryptography/selected-algorithms-2022>.

[3] V. Lyubashevsky, L. Ducas, E. Kiltz, T. Lepoint, P. Schwabe, G. Seiler, D. Stehlé, and S. Bai, “CRYSTALS-DILITHIUM,” National Institute of Standards and Technology, Tech. Rep., 2022, available at <https://csrc.nist.gov/Projects/post-quantum-cryptography/selected-algorithms-2022>.

[4] T. Prest, P.-A. Fouque, J. Hoffstein, P. Kirchner, V. Lyubashevsky, T. Pornin, T. Ricosset, G. Seiler, W. Whyte, and Z. Zhang, “FALCON,” National Institute of Standards and Technology, Tech. Rep., 2022, available at <https://csrc.nist.gov/Projects/post-quantum-cryptography/selected-algorithms-2022>.

[5] P. C. Kocher, J. Jaffe, and B. Jun, “Differential power analysis,” in *CRYPTO’99*, ser. LNCS, M. J. Wiener, Ed., vol. 1666. Springer, Heidelberg, Aug. 1999, pp. 388–397.

[6] E. Brier, C. Clavier, and F. Olivier, “Correlation power analysis with a leakage model,” in *CHES 2004*, ser. LNCS, M. Joye and J.-J. Quisquater, Eds., vol. 3156. Springer, Heidelberg, Aug. 2004, pp. 16–29.

[7] P. L. Montgomery, “Modular multiplication without trial division,” *Mathematics of computation*, vol. 44, no. 170, pp. 519–521, 1985.

[8] P. Barrett, “Implementing the Rivest Shamir and Adleman public key encryption algorithm on a standard digital signal processor,” in *CRYPTO’86*, ser. LNCS, A. M. Odlyzko, Ed., vol. 263. Springer, Heidelberg, Aug. 1987, pp. 311–323.

[9] A. Abdulrahman, V. Hwang, M. J. Kannwischer, and A. Sprenkels, “Faster kyber and dilithium on the cortex-M4,” in *ACNS 22*, ser. LNCS, G. Ateniese and D. Venturi, Eds., vol. 13269. Springer, Heidelberg, Jun. 2022, pp. 853–871.

[10] D. O. C. Greconici, M. J. Kannwischer, and A. Sprenkels, “Compact dilithium implementations on cortex-M3 and cortex-M4,” *IACR TCHES*, vol. 2021, no. 1, pp. 1–24, 2021, <https://tches.iacr.org/index.php/TCHES/article/view/8725>.

[11] E. Alkim, Y. A. Bilgin, M. Cenk, and F. Gérard, “Cortex-m4 optimizations for  $\{R, M\}$  lwe schemes,” *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 336–357, 2020.

[12] L. Botros, M. J. Kannwischer, and P. Schwabe, “Memory-efficient high-speed implementation of Kyber on cortex-M4,” in *AFRICACRYPT 19*, ser. LNCS, J. Buchmann, A. Nitaj, and T. eddine Rachidi, Eds., vol. 11627. Springer, Heidelberg, Jul. 2019, pp. 209–228.

[13] G. Seiler, “Faster avx2 optimized ntt multiplication for ring-lwe lattice cryptography,” *Cryptology ePrint Archive*, 2018.

[14] T. Plantard, “Efficient word size modular arithmetic,” *IEEE Transactions on Emerging Topics in Computing*, vol. 9, no. 3, pp. 1506–1518, 2021.

[15] J. Huang, J. Zhang, H. Zhao, Z. Liu, R. C. C. Cheung, Ç. K. Koç, and D. Chen, “Improved plantard arithmetic for lattice-based cryptography,” *IACR TCHES*, vol. 2022, no. 4, pp. 614–636, 2022.

[16] D. Heinz, M. J. Kannwischer, G. Land, T. Pöppelmann, P. Schwabe, and A. Sprenkels, “First-order masked kyber on ARM cortex-M4,” *Cryptology ePrint Archive*, Report 2022/058, 2022, <https://eprint.iacr.org/2022/058>.

[17] J. W. Bos, M. Gourjon, J. Renes, T. Schneider, and C. van Vredendaal, “Masking kyber: First- and higher-order implementations,” *IACR TCHES*, vol. 2021, no. 4, pp. 173–214, 2021, <https://tches.iacr.org/index.php/TCHES/article/view/9064>.

[18] T. Fritzzmann, M. V. Beirendonck, D. B. Roy, P. Karl, T. Schamberger, I. Verbauwhede, and G. Sigl, “Masked accelerators and instruction set extensions for post-quantum cryptography,” *IACR TCHES*, vol. 2022, no. 1, pp. 414–460, 2022.

[19] S. Özeren and O. Yayla, “Methods for masking crystals-kyber against side-channel attacks,” in *2023 16th International Conference on Information Security and Cryptology (ISCTürkiye)*. IEEE, 2023, pp. 1–6.

[20] V. Migliore, B. Gérard, M. Tibouchi, and P.-A. Fouque, “Masking dilithium: Efficient implementation and side-channel evaluation,” in *Applied Cryptography and Network Security: 17th International Conference, ACNS 2019, Bogota, Colombia, June 5–7, 2019, Proceedings 17*. Springer, 2019, pp. 344–362.

[21] M. Azouaoui, O. Bronchain, G. Cassiers, C. Hoffmann, Y. Kuzovkova, J. Renes, T. Schneider, M. Schönauer, F.-X. Standaert, and C. van Vredendaal, “Protecting dilithium against leakage revisited sensitivity analysis and improved implementations,” *IACR TCHES*, vol. 2023, no. 4, pp. 58–79, 2023.

[22] J.-S. Coron, F. Gérard, M. Trannoy, and R. Zeitoun, “Improved gadgets for the high-order masking of dilithium,” *IACR TCHES*, vol. 2023, no. 4, pp. 110–145, 2023.

[23] M. V. Beirendonck, J.-P. D’Anvers, A. Karmakar, J. Balasch, and I. Verbauwhede, “A side-channel-resistant implementation of saber,” *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, vol. 17, no. 2, pp. 1–26, 2021.

[24] S. Kundu, J.-P. D’Anvers, M. Van Beirendonck, A. Karmakar, and I. Verbauwhede, “Higher-order masked saber,” in *International Conference on Security and Cryptography for Networks*. Springer, 2022, pp. 93–116.

[25] O. Reparaz, S. S. Roy, R. de Clercq, F. Vercauteren, and I. Verbauwhede, “Masking ring-LWE,” *Journal of Cryptographic Engineering*, vol. 6, no. 2, pp. 139–153, Jun. 2016.

[26] O. Bronchain and G. Cassiers, “Bitslicing arithmetic/boolean masking conversions for fun and profit with application to lattice-based KEMs,” *IACR TCHES*, vol. 2022, no. 4, pp. 553–588, 2022.

[27] J.-S. Coron, F. Gérard, S. Montoya, and R. Zeitoun, “High-order polynomial comparison and masking lattice-based encryption,” *IACR TCHES*, vol. 2023, no. 1, pp. 153–192, 2023.

[28] D. Heinz and G. Dreo Rodosek, “Fast first-order masked ntru,” in *International Workshop on Constructive Side-Channel Analysis and Secure Design*. Springer, 2023, pp. 127–148.

[29] M. J. Kannwischer, J. Rijneveld, P. Schwabe, and K. Stoffelen, “pqm4: Testing and benchmarking NIST PQC on ARM cortex-M4,” *Cryptology ePrint Archive*, Report 2019/844, 2019, <https://eprint.iacr.org/2019/844>.

[30] R. Primas, P. Pessl, and S. Mangard, “Single-trace side-channel attacks on masked lattice-based encryption,” in *CHES 2017*, ser. LNCS, W. Fischer and N. Homma, Eds., vol. 10529. Springer, Heidelberg, Sep. 2017, pp. 513–533.

[31] I.-J. Kim, T.-H. Lee, J. Han, B.-Y. Sim, and D.-G. Han, “Novel single-trace ML profiling attacks on NIST 3 round candidate dilithium,” *Cryptology ePrint Archive*, Report 2020/1383, 2020, <https://eprint.iacr.org/2020/1383>.

[32] Z. Xu, O. Pemberton, S. S. Roy, D. Oswald, W. Yao, and Z. Zheng, “Magnifying side-channel leakage of lattice-based cryptosystems with chosen ciphertexts: The case study of kyber,” *IEEE Transactions on Computers*, vol. 71, no. 9, pp. 2163–2176, 2021.

[33] O. Bronchain, M. Azouaoui, M. ElGhamrawy, J. Renes, and T. Schneider, “Exploiting small-norm polynomial multiplication with physical attacks: Application to crystals-dilithium,” *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2024, no. 2, pp. 359–383, 2024.

[34] L. Backlund, K. Ngo, J. Gärtner, and E. Dubrova, “Secret key recovery attack on masked and shuffled implementations of crystals-kyber and saber,” in *International Conference on Applied Cryptography and Network Security*. Springer, 2023, pp. 159–177.

- [35] E. Dubrova, K. Ngo, J. Gärtner, and R. Wang, "Breaking a fifth-order masked implementation of crystals-kyber by copy-paste," in *Proceedings of the 10th ACM Asia Public-Key Cryptography Workshop*, 2023, pp. 10–20.
- [36] S. Marzougui, V. Ulitzsch, M. Tibouchi, and J.-P. Seifert, "Profiling side-channel attacks on dilithium: A small bit-fiddling leak breaks it all," *Cryptology ePrint Archive*, 2022.
- [37] E. Karabulut, E. Alkim, and A. Aysu, "Single-trace side-channel attacks on  $\omega$ -small polynomial sampling: With applications to NTRU, NTRU prime, and CRYSTALS-DILITHIUM," *Cryptology ePrint Archive*, Report 2022/494, 2022, <https://eprint.iacr.org/2022/494>.
- [38] F. Aydin, A. Aysu, M. Tiwari, A. Gerstlauer, and M. Orshansky, "Horizontal side-channel vulnerabilities of post-quantum key exchange and encapsulation protocols," *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 20, no. 6, pp. 1–22, 2021.
- [39] C. Mujdei, L. Wouters, A. Karmakar, A. Beckers, J. M. Bermudo Mera, and I. Verbauwhede, "Side-channel analysis of lattice-based post-quantum cryptography: Exploiting polynomial multiplication," *ACM Transactions on Embedded Computing Systems*, vol. 23, no. 2, pp. 1–23, 2024.
- [40] Z. Chen, E. Karabulut, A. Aysu, Y. Ma, and J. Jing, "An efficient non-profiled side-channel attack on the crystals-dilithium post-quantum signature," in *2021 IEEE 39th International Conference on Computer Design (ICCD)*. IEEE, 2021, pp. 583–590.
- [41] T. Tosun and E. Savas, "Zero-value filtering for accelerating non-profiled side-channel attack on incomplete ntt based implementations of lattice-based cryptography," *IEEE Transactions on Information Forensics and Security*, 2024.
- [42] H. Steffen, G. Land, L. Kogelheide, and T. Güneysu, "Breaking and protecting the crystal: Side-channel analysis of Dilithium in hardware," *Cryptology ePrint Archive*, Report 2022/1410, 2022, <https://eprint.iacr.org/2022/1410>.
- [43] M. J. Kannwischer, J. Rijneveld, P. Schwabe, and K. Stoffelen, "pqm4: Testing and benchmarking nist pqc on arm cortex-m4," 2019.
- [44] N. Bindel, S. Akleylek, E. Alkim, P. S. L. M. Barreto, J. Buchmann, E. Eaton, G. Gutoski, J. Kramer, P. Longa, H. Polat, J. E. Ricardini, and G. Zanon, "qTESLA," National Institute of Standards and Technology, Tech. Rep., 2019, available at <https://csrc.nist.gov/projects/post-quantum-cryptography/post-quantum-cryptography-standardization/round-2-submissions>.
- [45] C. Chen, O. Danba, J. Hoffstein, A. Hulsing, J. Rijneveld, J. M. Schanck, P. Schwabe, W. Whyte, Z. Zhang, T. Saito, T. Yamakawa, and K. Xagawa, "NTRU," National Institute of Standards and Technology, Tech. Rep., 2020, available at <https://csrc.nist.gov/projects/post-quantum-cryptography/post-quantum-cryptography-standardization/round-3-submissions>.
- [46] D. J. Bernstein, B. B. Brumley, M.-S. Chen, C. Chuengsatiansup, T. Lange, A. Marotzke, B.-Y. Peng, N. Tuveri, C. van Vredendaal, and B.-Y. Yang, "NTRU Prime," National Institute of Standards and Technology, Tech. Rep., 2020, available at <https://csrc.nist.gov/projects/post-quantum-cryptography/post-quantum-cryptography-standardization/round-3-submissions>.
- [47] V. Lyubashevsky, C. Peikert, and O. Regev, "On ideal lattices and learning with errors over rings," in *EUROCRYPT 2010*, ser. LNCS, H. Gilbert, Ed., vol. 6110. Springer, Heidelberg, May / Jun. 2010, pp. 1–23.
- [48] O. Regev, "On lattices, learning with errors, random linear codes, and cryptography," in *37th ACM STOC*, H. N. Gabow and R. Fagin, Eds. ACM Press, May 2005, pp. 84–93.
- [49] V. Lyubashevsky, "Fiat-Shamir with aborts: Applications to lattice and factoring-based signatures," in *ASIACRYPT 2009*, ser. LNCS, M. Matsui, Ed., vol. 5912. Springer, Heidelberg, Dec. 2009, pp. 598–616.
- [50] H. Steffen, G. Land, L. Kogelheide, and T. Güneysu, "Breaking and protecting the crystal: Side-channel analysis of dilithium in hardware," in *International Conference on Post-Quantum Cryptography*. Springer, 2023, pp. 688–711.
- [51] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex fourier series," *Mathematics of computation*, vol. 19, no. 90, pp. 297–301, 1965.
- [52] W. M. Gentleman and G. Sande, "Fast fourier transforms: for fun and profit," in *Proceedings of the November 7-10, 1966, fall joint computer conference*, 1966, pp. 563–578.
- [53] V. Lyubashevsky and G. Seiler, "Ntru: Truly fast ntru using ntt," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 180–201, 2019.
- [54] C.-M. M. Chung, V. Hwang, M. J. Kannwischer, G. Seiler, C.-J. Shih, and B.-Y. Yang, "Ntt multiplication for ntt-unfriendly rings: New speed records for saber and ntru on cortex-m4 and avx2," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 159–188, 2021.
- [55] A. Abdulrahman, J.-P. Chen, Y.-J. Chen, V. Hwang, M. J. Kannwischer, and B.-Y. Yang, "Multi-moduli ntt for saber on cortex-m3 and cortex-m4," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 127–151, 2022.
- [56] E. Prouff, M. Rivain, and R. Bevan, "Statistical analysis of second order differential power analysis," *IEEE Transactions on computers*, vol. 58, no. 6, pp. 799–811, 2009.
- [57] F.-X. Standaert, N. Veyrat-Charvillon, E. Oswald, B. Gierlichs, M. Medwed, M. Kasper, and S. Mangard, "The world is not enough: Another look on second-order DPA," in *ASIACRYPT 2010*, ser. LNCS, M. Abe, Ed., vol. 6477. Springer, Heidelberg, Dec. 2010, pp. 112–129.
- [58] M. Joye, P. Paillier, and B. Schoenmakers, "On second-order differential power analysis," in *Cryptographic Hardware and Embedded Systems—CHES 2005: 7th International Workshop, Edinburgh, UK, August 29–September 1, 2005. Proceedings 7*. Springer, 2005, pp. 293–308.
- [59] A. Aikata, A. Basso, G. Cassiers, A. C. Mert, and S. S. Roy, "Kavach: Lightweight masking techniques for polynomial arithmetic in lattice-based cryptography," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 366–390, 2023.
- [60] B. Gierlichs, L. Batina, P. Tuyls, and B. Preneel, "Mutual information analysis: A generic side-channel distinguisher," in *International Workshop on Cryptographic Hardware and Embedded Systems*. Springer, 2008, pp. 426–442.
- [61] Y.-T. Kuo and A. Takayasu, "A lattice attack on crystals-kyber with correlation power analysis," in *International Conference on Information Security and Cryptology*. Springer, 2023, pp. 202–220.

## APPENDIX. ABSOLUTE DIFFERENCE COMBINATION FUNCTION

Absolute value combination function [58] is defined as follows.

$$C_{abs}(\mathcal{L}(X^0), \mathcal{L}(X^1)) = |\mathcal{L}(X^0) - \mathcal{L}(X^1)| \quad (19)$$

In the following analysis including Figure 12, Figure 13, Table 6, Table 7, and Table 8, for the sake of simplicity we took  $\mu_0 = \mu_1 = 0$  regarding the noise terms in  $\mathcal{L}(X^0)$  and  $\mathcal{L}(X^1)$  (see Equation (4)). Note that, our results hold as long as  $\mu_0 = \mu_1$ . Estimation of the corresponding optimal prediction function – so-called  $\hat{f}_{opt}(i)$  – is performed with 100 k samples uniformly taken for  $X^0, X^1$ , and the noise taken from  $\mathcal{N}(0, \sigma)$  for each  $i$ . Results that depend on the  $\hat{f}_{opt}$  for  $q = 8380417$  are not presented due to the computational complexity of experiments. However, they can be anticipated from the results for the other studied primes presented in this section and the ones corresponding to the mean-free product (Figure 3 and Table 5).

### A. OPTIMAL CORRELATION FOR $C_{ABS}$

Demonstrated in Figure 12.

### B. ACCURACY OF ABSOLUTE VALUE PREDICTION FUNCTION FOR $C_{ABS}$

It can be seen in Figure 13 that  $\hat{f}_{opt}$  for  $C_{abs}$ , is an affine function of  $f_{abs}$  with some noise. Table 6, Table 7, and Table 8 present the estimations for  $\hat{\rho}(f_{abs}(X), \hat{f}_{opt}(X))$  in this configuration. Observe that  $|\hat{\rho}| > 0.99$  for the evaluated  $q$  and  $\beta$  couples mentioned in Section IV, allowing us to conclude that  $f_{abs}$  can be effectively used with  $C_{abs}$ .

Additionally, we provide the correlation results for  $q = 8380417$  in Figure 14.

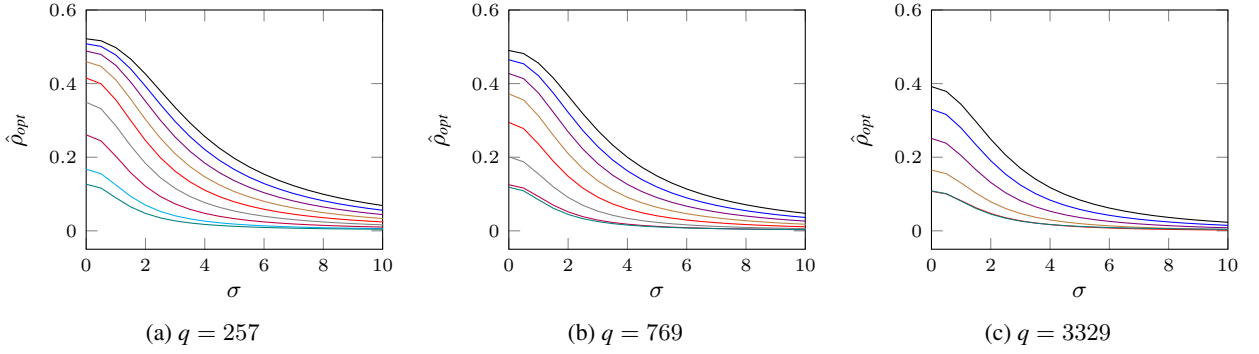


FIGURE 12: Optimal correlation for  $\mathcal{C}_{abs}$  with respect to the noise standard deviation  $\sigma$  for different  $q$  and  $\beta$  and reduction algorithms. Estimations are performed with 1 million samples uniformly taken for  $X^0$  and  $X^1$ .

- ★ Reduction to  $[-q/2, q/2]$  for  $q = 257$ ,  $q = 769$  and  $q = 3329$ :  $\beta = 16$  (black),  $\beta = 15$  (blue),  $\beta = 14$  (violet),  $\beta = 13$  (brown),  $\beta = 12$  (red),  $\beta = 11$  (gray),  $\beta = 10$  (purple),  $\beta = 9$  (cyan)
- ★ Reduction to  $[0, q)$  (teal)

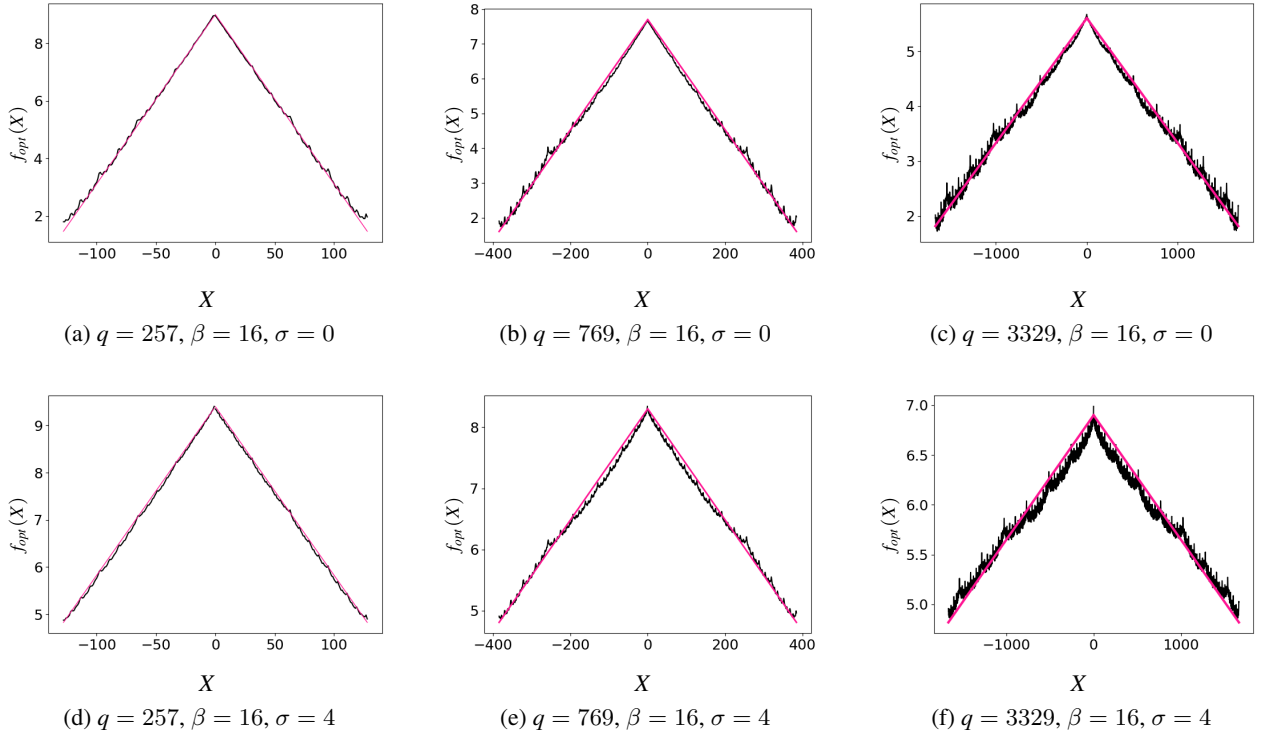


FIGURE 13: Visualization of optimal prediction function  $\hat{f}_{opt}$  for the central reduction range  $[-q/2, q/2]$  with respect to absolute difference combination function.  $\hat{f}_{opt}(X)$  (black),  $c'_0 \cdot |X| + c'_1$  (pink) for some  $c'_0$  and  $c'_1$ .

$\beta \backslash \sigma$	0	2	4	6	8	10
<b>16</b>	<b>-0.999</b>	<b>-0.999</b>	<b>-0.999</b>	<b>-0.999</b>	<b>-0.999</b>	<b>-0.999</b>
15	-0.999	-0.999	-0.999	-0.999	-0.999	-0.998
14	-0.998	-0.998	-0.998	-0.998	-0.998	-0.998
13	-0.998	-0.998	-0.997	-0.997	-0.996	-0.996
12	-0.996	-0.996	-0.994	-0.993	-0.993	-0.992
11	-0.992	-0.990	-0.986	-0.984	-0.983	-0.982
10	-0.975	-0.964	-0.957	-0.952	-0.949	-0.947
9	-0.864	-0.834	-0.818	-0.809	-0.801	-0.793

TABLE 6: Estimations of  $\hat{\rho}(f_{abs}(X), \hat{f}_{opt}(X))$  for  $\mathcal{C}_{abs}$ ,  $q = 257$  and different  $\beta$  and  $\sigma$ . The estimations are performed with 1 million uniformly random samples for  $X$  while the reduction is central.

$\beta \backslash \sigma$	0	2	4	6	8	10
<b>16</b>	<b>-0.998</b>	<b>-0.998</b>	<b>-0.998</b>	<b>-0.998</b>	<b>-0.998</b>	<b>-0.998</b>
15	-0.997	-0.998	-0.997	-0.997	-0.996	-0.996
14	-0.996	-0.996	-0.995	-0.994	-0.994	-0.994
13	-0.993	-0.992	-0.990	-0.989	-0.988	-0.988
12	-0.985	-0.979	-0.974	-0.972	-0.971	-0.970
11	-0.941	-0.923	-0.912	-0.908	-0.905	-0.902
10	-0.673	-0.645	-0.632	-0.626	-0.620	-0.614

TABLE 7: Estimations of  $\hat{\rho}(f_{abs}(X), \hat{f}_{opt}(X))$  for  $\mathcal{C}_{abs}$ ,  $q = 769$  and different  $\beta$  and  $\sigma$ . The estimations are performed with 1 million uniformly random samples for  $X$  while the reduction is central.

$\beta \backslash \sigma$	0	2	4	6	8	10
<b>16</b>	<b>-0.996</b>	<b>-0.995</b>	<b>-0.993</b>	<b>-0.993</b>	<b>-0.992</b>	<b>-0.992</b>
15	-0.992	-0.989	-0.986	-0.985	-0.984	-0.983
14	-0.978	-0.971	-0.965	-0.962	-0.961	-0.959
13	-0.913	-0.893	-0.882	-0.877	-0.872	-0.869
12	-0.562	-0.551	-0.544	-0.539	-0.534	-0.529

TABLE 8: Estimations of  $\hat{\rho}(f_{abs}(X), \hat{f}_{opt}(X))$  for  $\mathcal{C}_{abs}$ ,  $q = 3329$  and different  $\beta$  and  $\sigma$ . The estimations are performed with 1 million uniformly random samples for  $X$  while the reduction is central.

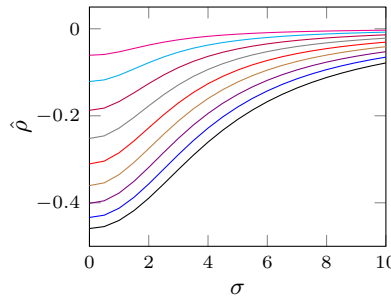


FIGURE 14: Estimations of  $\hat{\rho}(f_{abs}(X), \mathcal{C}_{abs}(X^0, X^1))$  with respect to the noise standard deviation  $\sigma$  for  $q = 8380417$  and different  $\beta$ . The estimations are performed with 1 million samples uniformly taken for  $X^0$  and  $X^1$ .

- Reduction to  $[-q/2, q/2]$  for  $q = 8380417$ :  $\beta = 32$  (black),  $\beta = 31$  (blue),  $\beta = 30$  (violet),  $\beta = 29$  (brown),  $\beta = 28$  (red),  $\beta = 27$  (gray),  $\beta = 26$  (purple),  $\beta = 25$  (cyan),  $\beta = 24$  (magenta).