

Zero-Value Filtering for Accelerating Non-Profiled Side-Channel Attack on Incomplete NTT based Implementations of Lattice-based Cryptography

Tolun Tosun^{1,2}, Erkay Savas¹

¹*Sabanci University Computer Science and Engineering* ²*Analog Devices*

Abstract—Lattice-based cryptographic schemes such as Crystals-Kyber and Dilithium are post-quantum algorithms selected to be standardized by NIST as they are considered to be secure against quantum computing attacks. The multiplication in polynomial rings is the most time-consuming operation in many lattice-based cryptographic schemes, which is also subject to side-channel attacks. While NTT-based polynomial multiplication is almost a norm in a wide range of implementations, a relatively new method, incomplete NTT is preferred to accelerate lattice-based cryptography, especially on some computing platforms that feature special instructions. In this paper, we present a novel, efficient and non-profiled power/EM side-channel attack targeting polynomial multiplication based on the incomplete NTT algorithm. We apply the attack on the Crystals-Dilithium signature algorithm and Crystals-Kyber KEM. We demonstrate that the method accelerates attack run-time when compared to the existing approaches. While a conventional non-profiled side-channel attack tests a much larger hypothesis set because it needs to predict two coefficients of secret polynomials together, we propose a much faster *zero-value filtering attack (ZV-FA)*, which reduces the size of the hypothesis set by targeting the coefficients individually. We also propose an effective and efficient validation and correction technique employing the inverse NTT to estimate and modify the mispredicted coefficients. Our experimental results show that we can achieve a speed-up of $1915\times$ over brute-force.

Index Terms—post-quantum cryptography; side-channel attack; correlation power analysis; multivariate mutual information analysis; crystals dilithium; crystals kyber;

I. INTRODUCTION

SECURITY of public-key cryptosystems relies on the hardness of well-known mathematical problems such as the discrete logarithm problem for the elliptic curve cryptography (ECC) [1], [2] and the digital signature algorithm (DSA) [3] or the integer factorization problem for RSA [4]. While those hard problems are conjectured to be secure against known cryptanalytic algorithms running on classical computers, it has been shown that Shor's algorithm [5] can solve them in polynomial time on a large-scale quantum computer.

To address the quantum threat, the National Institute of Standards and Technology (NIST) has announced the standardization process for post-quantum public-key cryptographic algorithms (PQC) in 2016. The standardization process covers quantum-resistant digital signature schemes, and public-key encryption and key-establishment algorithms. Currently, the contest is at the fourth round with already standardized algorithms. Lattice-based schemes, based on various hard lattice problems, facilitate the construction of quantum resilient

public-key cryptography with a promising level of efficiency. Among the winners, the lattice-based digital signature algorithm Crystals-Dilithium [6] is based on the Module-LWE [7] and Module-SIS (MSIS) problems, while Crystals-Kyber [8] is MLWE based key encapsulation mechanism (KEM). As Kyber and Dilithium are members of the same family, Crystals, they have several building blocks in common.

In cryptoanalysis, side-channel attacks (SCA) are the ones that target the weaknesses in implementations rather than algorithm specifications, by collecting side information such as running time or power consumption that can leak sensitive (intermediate) information during the execution of the targeted cryptographic operation. Side-channel attacks are considered as one of the main threats, particularly for embedded devices because of the simplicity of side-information collection, such as IoT chips, which sign sensor data before transmission. The Correlation Power Analysis (CPA) is proposed in [9], which models the power consumption of the device under test and measures the correlation of the model with real-world data to test secret value hypotheses. The power leakage of the device/implementation is often modeled with the Hamming Weight (HW)/Hamming Distance (HD) of/between the intermediate data. The Mutual Information Analysis (MIA) [10] is another efficient side-channel *distinguisher*, based on information theory and Shannon entropy. The Electromagnetic (EM) side-channels [11] are similar to power analysis as any attack suit designed for power leakage can be practiced with EM leakage while it can supply more precise information about the sensitive intermediate data. Masking is one of the most promising countermeasures against power/EM attacks, which randomizes the intermediate data with secret sharing so that characteristics of sensitive data are not reflected in power consumption.

JIL Rating [12] is a widely used metric to assess the complexity of side-channel attacks; the higher the JIL score, the harder to perform the attack. As the time needed to apply the attack is a factor in the overall rating, both the number of traces and the attack run-time affect the rating of the attack, constituting an important motivation for this work.

Needless to say, the side-channel security of post-quantum public cryptography is essential as well since post-quantum algorithms are intended to replace the existing public-key standards soon and the usage of public-key cryptography in embedded devices will be potentially more extensive. For example, a secure firmware update on an embedded device

relies on the security of the employed digital signature while embedded devices are open to timing, and power/EM attacks by nature. With increasing interest, several attacks and countermeasures have been proposed for PQC candidates in the literature.

Polynomial multiplication is the core operation for practical constructions of lattice-based cryptography, which are based on ring-learning with errors (R-LWE) problem [13]. Most implementations utilize number theoretic transform (NTT) for efficient polynomial arithmetic [14]. A technique referred as *incomplete NTT* is introduced to handle rings of special structures as well as for efficiency [15]–[17] in implementations of various lattice-based cryptography algorithms. Our study targets the incomplete NTT operation specifically.

Table I summarizes the related side-channel attacks against lattice-based schemes from the literature. Among the attack types, the profiled class forms the majority, where we require a device identical to the one targeted by the attacker, who tries to characterize the leakage when executing a cryptographic algorithm with a known secret key. In other words, the attacker needs to have more capability in a profiled attack compared to the non-profiled class. Machine learning-based approaches are gaining popularity in the design of profiled attacks for lattice-based cryptography [18], [19]. In addition, Primas et al. [20] present a notable study by combining the side-channel leakage of NTT computation with the belief propagation algorithm to conduct a single-trace profiled attack.

As for the non-profiled class, the polynomial multiplication is the most attractive target operation [21]–[23]. Steffen et al. [23] conduct an attack on a hardware implementation of Dilithium. Chen et al. [22] target the reference implementation of Dilithium, concentrating on improving the runtime performance of non-profiled attack, as the conventional approach requires brute-force effort over 23-bit secrets based on the coefficient modulus length of Dilithium. Mujdei et al. [21] attack ARM M4 implementation of Kyber [17], which has a very similar NTT implementation with Dilithium, based on [16]. The authors of [21] show that the secret coefficients must be predicted in pairs since the incomplete NTT algorithm is used in the polynomial multiplication of the targeted implementation.

In our study, we present a more efficient non-profiled attack on the incomplete NTT implementation, which facilitates that the coefficients of the secret polynomials can be predicted individually. To show its efficacy, we use the very recent and fast implementation of Dilithium and Kyber on ARM M4 [17] as in [21]. We present several non-profiled side-channel attacks targeting the multiplication in the incomplete NTT domain and finally develop a much more efficient approach exploiting the zero-valued coefficients of the known operand of the incomplete NTT multiplication, with application to Dilithium and Kyber schemes.

A. Main Contributions

We can list our contributions as follows:

We present a novel non-profiled power/EM attack against incomplete NTT-based implementations of polynomial

multiplication in Lattice-based Cryptography, referred to as Zero-Value Filtering Attack (ZV-FA). Our approach is efficient as it decreases the number of hypotheses significantly by introducing a filtering technique based on zero-value coefficients in the known input/output polynomials of the operation targeted by the side-channel attack.

We present an efficient validation technique for estimating and correcting mispredicted values for attacking secret polynomials with short coefficients. The method not only ensures full accuracy on the estimated secret polynomials but also accelerates the attack run time by trading off the number of traces.

We show that using short secret key polynomials in lattice-based cryptography can be exploited to accelerate side-channel attacks.

We implement the ZV-FA with the validation technique on the pqm4 [17] implementations of Dilithium and Kyber [16]. Our experiments demonstrate that a moderate increase in the number of traces can decrease the attack run-time significantly. It is experimentally shown through EM side-channel that, a speed-up of up to three orders of magnitude in attack run-time can be achieved over a conventional CPA targeting the polynomial multiplication.

We experimentally show that our approach is also favorable in the presence of masking, by applying ZV-FA to a protected implementation of Kyber.

II. NOTATION

Matrices are represented by bold uppercase letters, such as \mathbf{A} , while vectors are represented by bold lowercase letters, such as \mathbf{b} . Sets are denoted by uppercase calligraphic letters, such as \mathcal{A} . Polynomials are denoted by lowercase italic letters, such as f . Depending on the context, polynomials may be represented together with their indeterminate, such as $f(x)$. Subscripts together with square brackets are used to denote element(s) of matrices and vectors, such as $\mathbf{A}_{[i:j]}$ and $\mathbf{s}_{[i]}$; elements of sets, such as $A_{[i]}$; coefficients of polynomials, such as $f_{[i]}$. The notation $\mathbf{A}_{[:,j]}$ denotes the j -th column vector of \mathbf{A} . Similarly, $A_{[:,i]}$ denotes the first i elements of the set \mathcal{A} .

Modular reductions are performed in a centered manner. Specifically, given an integer i and a modulus q , the operation $i^{\flat} = i \pmod{q}$ maps i to a unique integer i^{\flat} in the range of $[-\lfloor bq/2c \rfloor; \lfloor bq/2c \rfloor]$ for odd q . The set of integers modulo q with centered reduction is denoted by \mathbb{Z}_q . On the other hand, the notation \mathbb{Z}_q^+ is used to denote the set of integers modulo q using the positive range, namely $[0; q-1]$. The ring of polynomials $\mathbb{Z}_q[x] = (X^n + 1)$, where elements are polynomials of the maximum degree of $n-1$, whose coefficients are modulo- q reduced, is denoted by R_q .

The dimensionality of matrices and vectors is shown in the superscript. For example, $R_q^{k \times l}$ represents a matrix of dimensions $k \times l$, whose elements are in R_q . Similarly, superscript is used for the matrices and vectors to express their dimensionality, such as $\mathbf{A}^{N \times M}$ and \mathbf{b}^N . The cardinality of the set \mathcal{A} is denoted by $|\mathcal{A}|$. The matrix-vector multiplication of the operands \mathbf{A} and \mathbf{b} is denoted by $\mathbf{A}\mathbf{b}$. Similarly, multiplication of the vector \mathbf{b} with scalar a is denoted by $a\mathbf{b}$.

Attack	Class	Algorithm	Implementation	Target Operation
this work	Non-profiled	Dilithium / Kyber	ARM M4 (masked)	polynomial multiplication
[22]	Non-profiled	Dilithium	Reference C	polynomial multiplication
[23]	Non-Profiled	Dilithium	HW	polynomial multiplication
[21]	Non-profiled	Kyber	ARM M4	polynomial multiplication
[18]	Profiled	Dilithium	Reference C	NTT
[19]	Profiled	Dilithium	Reference C	bit-unpacking
[23]	Profiled	Dilithium	HW	decoding / NTT
[24]	Profiled	Dilithium	Reference C / ARM M4	small polynomial sampling
[25]	Profiled	Dilithium	Reference C	decompose
[26]	Profiled	Kyber	Reference C / ARM M4	NTT
[20]	Profiled	R-LWE Encryption	ARM M4 (masked) [27]	NTT

TABLE I: Related Side-Channel Attacks from the Literature

Polynomial multiplication is denoted by the standard symbol \cdot , while element-wise multiplication of two vectors is denoted by \odot . In some cases, the symbol \odot is used to represent integer multiplication to support the narrative. The symbol \times denotes the Cartesian product between sets, such as $A \times B$.

The notation (also referred to as infinity norm) $\|s\|_\infty$ is used to represent the maximum coefficient of the polynomial s in absolute value, whose elements are reduced in a centered manner. Similarly, $\|s\|_\infty$ is the maximum of the maximum absolute values of coefficients of the polynomials in the vector s . The set S consists of polynomials $w \in R_q$ with $\|w\|_\infty \leq \epsilon$, where ϵ is a (relatively small) positive integer, referred to as the set of *short polynomials*. B denotes the central binomial distribution over R_q , where $\|w\|_\infty \leq \epsilon$ as well for $w \in B$. Another subset of R_q is denoted by B , which consists of polynomials with exactly ϵ coefficients that are either -1 or 1, and the rest is zero. $\mathcal{F}_0; 1g^N$ denotes the set of N -bit strings. The operator \mathcal{F} denotes uniformly random sampling from the set on the right-hand side, such as $\mathcal{F}_0; 1g^N$. A prediction to a secret a is denoted by *wide hat*, such as \hat{a} .

III. SIDE-CHANNEL ATTACK OVERVIEW AND DISTINGUISHERS

Main steps of a non-profiled side-channel attack can be summarized as follows:

The attacker observes N cryptographic operations involving the secret key and records the power consumption of the victim device. M points are sampled in time at each observation. Power samples are stored in the matrix $\mathbf{T}^{N \times M}$ while \mathbf{p}^N is the vector of known variables.

A point of interest (PoI) is selected by the adversary. The PoI should be a function of a known variable that changes at each experiment and the attacked secret that remains the same for all experiments.

A set of predictions is prepared, denoted by K . Then, intermediate value matrix $\mathbf{V}^{N \times |K|}$ is computed w.r.t. each hypothesis; *i.e.*, $\mathbf{V}_{[i;j]}$ is the value of the PoI computed using $\mathbf{p}_{[i]}$ and $K_{[j]}$.

\mathbf{V} is mapped to the hypothetical power consumption matrix, $\mathbf{H}^{N \times |K|}$ by applying a power consumption model. In other words, $\mathbf{H}_{[i;j]} = \text{HW}(\mathbf{V}_{[i;j]})$, in case HW is chosen.

Each hypothesis $K_{[j]}$ is tested, *i.e.* scored, with the selected distinguisher, by considering the relationship

between $\mathbf{H}_{[i;j]}$ and $\mathbf{T}_{[i;j_1]}$ (and $\mathbf{T}_{[i;j_2]}$ for second-order case). Usually, the maximum score over all j_1 (and j_2) is assigned.

In side-channel literature, distinguishers are statistical tools that are used to rank predictions (*i.e.*, hypotheses) for secret keys, by exploiting the dependency between the data processed by a cryptographic implementation and its power consumption. The distinguishers employed in this study are presented in the following subsections.

A. Correlation Power Analysis (CPA)

CPA [9] is a widely-used side-channel distinguisher, based on Pearson's correlation. For random variables X and Y , the correlation coefficient is estimated by

$$\hat{c}(X; Y) = \frac{\text{cov}(X; Y)}{\hat{\text{std}}(X) \hat{\text{std}}(Y)}; \quad (1)$$

To utilize correlation as a side-channel distinguisher, $\hat{c}(\mathbf{H}_{[i;j]}; \mathbf{T}_{[i;j_1]})$ is computed.

B. Higher-Order CPA (HOCPA)

To perform a higher-order attack using CPA, multiple samples must be combined using a pre-processing function. The state-of-art pre-processing function is the normalized product [28], defined as follows for the second-order case

$$\mathbf{T}_{[i;j_1 j_2]}^0 = (\mathbf{T}_{[i;j_1]} \overline{\mathbf{T}_{[i;j_1]}}) (\mathbf{T}_{[i;j_2]} \overline{\mathbf{T}_{[i;j_2]}}) \text{ for } 0 < N \quad (2)$$

which combines the leakage at time samples j_1 and j_2 . Then, $\hat{c}(\mathbf{H}_{[i;j]}; \mathbf{T}_{[i;j_1 j_2]}^0)$ is computed.

C. Mutual Information Analysis (MIA)

MIA [10] is an information-theoretic side-channel distinguisher. Let $H(X)$ denote the entropy of a random variable X on a discrete space X , and let $H(X|Y)$ denote the conditional entropy for X and another random variable Y on a discrete space Y . $P(\cdot)$ denotes the probability function. The mutual information between X and Y is defined as

$$\begin{aligned} I(X; Y) &= H(X) - H(X|Y) \\ &= \sum_{x \in X} \sum_{y \in Y} P(X=x; Y=y) \log_2 \frac{P(X=x; Y=y)}{P(X=x)P(Y=y)} \end{aligned} \quad (3)$$

NIST Security Level		2	3	5
Parameter	Meaning			
τ	number of 1 in c	39	49	60
k, l	dimensions of \mathbf{A}	(4,4)	(6,5)	(8,7)
η	coefficient range of $\mathbf{s}_1, \mathbf{s}_2$	2	4	2
	Expected #repetitions	4.25	5.1	3.85

TABLE II: Dilithium parameter set

Verbally, the entropy in X not covered by Y , namely $H(X|Y)$, is subtracted from $H(X)$ to formulate the mutual information in between. Similar to Section III-A, $I(\mathbf{H}_{[:;l]}; \mathbf{T}_{[:;j_1]})$ is computed in side-channel analysis.

D. Multivariate Mutual Information Analysis (MMIA)

MMIA [29] is a generalization of MIA to high-order side-channel attacks. The multivariate mutual information between three random variables X , Y , and Z is given by

$$I(X; Y; Z) = I(Y; Z) - I(Y; Z|X) \quad (4)$$

where

$$I(Y; Z|X) = \sum_{x \in \mathcal{X}} P(X = x) I(Y; Z|X = x) \quad (5)$$

By plugging two power samples in time into the above equation, a second-order side-channel distinguisher is achieved, namely $I(\mathbf{H}_{[:;l]}; \mathbf{T}_{[:;j_1]}, \mathbf{T}_{[:;j_2]})$.

IV. LATTICE-BASED POST-QUANTUM CRYPTOGRAPHY

In this section, we present the lattice-based post-quantum cryptography algorithms, Dilithium and Kyber.

A. Dilithium

Crystals: Dilithium [6] is a lattice-based post-quantum digital signature scheme based on the hardness of MLWE and MSIS problems. It operates over the ring of polynomials R_q with dimension $n = 256$ and $q = 8380417 = 2^{23} \cdot 2^{13} + 1$. The rest of the parameter set used by Dilithium can be found in Table II. Dilithium's key generation creates an MLWE instance by the equation $\mathbf{t} := \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2$. The matrix of polynomials $\mathbf{A} \in R_q^{k \times l}$ is part of the public key as well as the vector of polynomials $\mathbf{t} \in R_q^k$ (in the full scheme, lower bits of coefficients of polynomials in \mathbf{t} are kept secret). The vectors of short polynomials $\mathbf{s}_1 \in S^l$, $\mathbf{s}_2 \in S^k$ forms the secret key.

The template of Dilithium's signature generation, given in Algorithm 1, applies the rejection sampling idea. Most of the signature procedure is implemented in a loop, which iterates until a valid and secure signature is found. The parameters are chosen such that the expected number of repetitions is small, as presented in Table II. Inside the signature loop, a challenge polynomial $c \in B$ is generated. The candidate signature $\mathbf{z} = \mathbf{y} + c\mathbf{s}_1$ is rejected and restarted in case the security and correctness checks at line 7 fail. We would like to note that, there is no difference between the template and the standard Dilithium from the perspective of this work.

Algorithm 1 Dilithium.Sign(sk, M)

```

1:  $\mathbf{z} := ?$ 
2: while  $\mathbf{z} = ?$  do
3:    $\mathbf{y} \in \mathcal{S}'_1$ 
4:    $\mathbf{w}_1 := \text{HighBits}(\mathbf{A}\mathbf{y}; 2^{-2})$ 
5:    $c \in B := H(\text{jj } \mathbf{w}_1)$ 
6:    $\mathbf{z} := \mathbf{y} + c\mathbf{s}_1$ 
7:   if  $\text{jjzjj}_1 > \tau$  or  $\text{jjLowBits}(\mathbf{A}\mathbf{y} - c\mathbf{s}_2)\text{jj}_1 > \tau$  then
8:      $(\mathbf{z}) := ?$ 
9: return  $(\mathbf{z}; c)$ 

```

NIST Security Level		1	3	5
Parameter	Meaning			
k	dimensions of \mathbf{A}	2	3	4
η	Parameter of CBD B	2	4	2

TABLE III: Kyber parameter set

B. Kyber

Crystals: Kyber [6] is a lattice-based post-quantum KEM signature scheme based whose security relies on the computational difficulty of the MLWE problem. Kyber's ring dimension $n = 256$ and the coefficient modulus $q = 3329 = 2^{11} + 2^{10} + 2^8 + 1$ for the operated ring of polynomials R_q . The rest of the related parameter set can be found in Table III. Similar to Dilithium, the public-private key pair of Kyber is generated by the MLWE equation $\mathbf{t} := \mathbf{A}\mathbf{s} + \mathbf{e}$, where $\mathbf{s} \in R_q^k$ is the secret key, $\mathbf{t} \in R_q^k$ and $\mathbf{A} \in R_q^{k \times k}$ forms the public key and $\mathbf{e} \in R_q^k$ is the noise vector. \mathbf{s} and \mathbf{e} are short polynomials, whose coefficients are sampled from the central binomial distribution B .

Algorithm 2 Kyber.CPAPKE.Dec($\mathbf{s}, ct=(\mathbf{u}, v)$)

```

1: return  $\text{Compress}_q(v - \mathbf{s}^T \mathbf{u}; 1)$ 

```

Algorithm 2 presents a simplified version of Kyber's key decapsulation. The ciphertext ct is a composition of two parts; $\mathbf{u} \in R_q^k$ and $v \in R_q$. The function Compress_q maps from R_q to $f_0; 1g^n$.

V. NUMBER THEORETIC TRANSFORM (NTT)

Number Theoretic Transform (NTT) is a special form of Fast-Fourier Transform (FFT) that operates over \mathbb{Z}_q instead of complex numbers \mathbb{C} . NTT allows efficient multiplication of polynomials over R_q . Representing a polynomial $a(x) \in R_q$ in the NTT domain can be viewed as an application of Chinese Remainder Theorem (CRT). Polynomial multiplication is achieved by element-wise multiplying the NTT representations of the operands:

$$a(x) \cdot b(x) = \text{NTT}^{-1}(\text{NTT}(a(x)) \cdot \text{NTT}(b(x))) \quad (6)$$

Most lattice-based crypto systems, including Dilithium, operates over the ring R_q . NTT in R_q requires $q \equiv 1 \pmod{2n}$, which ensures a primitive $2n$ -th root of unity ω_{2n} exists in \mathbb{Z}_q for which $\omega_{2n}^n = -1 \pmod{q}$, referred to as *negacyclic NTT*. In case n is a power of 2, NTT can be computed efficiently

by splitting the polynomial to half of its size in recursive manner until linear degree is reached. The transformation at each layer can be efficiently implemented using *Cooley-Tukey* (CT) butterfly circuit [30]. For degree- $n=2^i$ polynomial $a(x) = a_0(x) + a_1(x) x^{n=2^{i+1}}$, the CT butterfly is defined by the map

$$a_0(x) + a_1(x) x^{n=2^{i+1}} \rightarrow (a_0(x) \quad a_1(x); a_0(x) + a_1(x)):$$

where n is an odd power of 2^n , called the *twiddle factor*. In this manner, full NTT computation requires $\log n$ layers. Most applications use *Gentleman-Sande* (GS) butterfly for computing the inverse NTT [31], although it is not necessarily required. Using CT or GS when n is a power of 2, computing NTT / inverse NTT takes $(n \log n)$ steps.

A. Incomplete NTT

For efficiency NTT can be computed for $m < \log n$ layers so polynomials are recursively splitted to degree- $n=2^m$ polynomial components, denoted by $\text{NTT}_m(a(x))$. The prerequisite for NTT_m is to have $q \equiv 1 \pmod{\frac{2^{\log n - m - 1}}{n}}$ for the negacyclic NTT [32].

Let $\mathbf{a} = \text{NTT}_m(a(x))$ and $\mathbf{b} = \text{NTT}_m(b(x))$. \mathbf{a} and \mathbf{b} are 2^m -dimensional vectors and $\mathbf{a}_{[i:]}, \mathbf{b}_{[i:]} \in \mathbb{Z}_q[x] = (x^{n=2^m})$ for $i < 2^m$, where n is some power of 2^n . Then, $a(x) b(x)$ can be computed through $\mathbf{a} \cdot \mathbf{b}$ as in Equation 6. In this case, element-wise multiplication refers to the multiplication of polynomials of degree $n=2^m - 1$, which is mostly implemented by the school-book algorithm, and there are 2^m such multiplications. For instance, when $m = \log(n=2)$, $n=2$ multiplications of degree-1 polynomials is performed, as presented in Algorithm 3.

Algorithm 3 Incomplete NTT Multiplication($\mathbf{s}^\ell, \mathbf{c}^\ell$)

```

1: for  $i = 0$  until  $n=2$  do
2:    $\mathbf{r}_{[i:0]} = \mathbf{s}_{[i:0]}^\ell \mathbf{c}_{[i:0]}^\ell + \mathbf{s}_{[i:1]}^\ell \mathbf{c}_{[i:1]}^\ell \quad i \bmod q$ 
3:    $\mathbf{r}_{[i:1]} = \mathbf{s}_{[i:0]}^\ell \mathbf{c}_{[i:1]}^\ell + \mathbf{s}_{[i:1]}^\ell \mathbf{c}_{[i:0]}^\ell \bmod q$ 
4: return  $\mathbf{r}$ 

```

VI. TARGET IMPLEMENTATION AND POINT OF INTEREST (POI)

We focus on the side-channel attack against incomplete NTT-based polynomial multiplication (Algorithm 3). One of the inputs is secret, and the attacker wants to learn it through a side-channel attack while one of the inputs is public. This case is directly applicable to Kyber and Dilithium, regarding the operations $\mathbf{s}^T \mathbf{u}$ and $\mathcal{C}\mathbf{s}_1$ (and $\mathcal{C}\mathbf{s}_2$), respectively. Needless to say, the attack aims to retrieve \mathbf{s} for Kyber and $\mathcal{C}\mathbf{s}_1$ and $\mathcal{C}\mathbf{s}_2$ for Dilithium. To simplify notation, we denote the attacked polynomial and its NTT representation by s^ℓ and \mathbf{s}^ℓ , respectively. Note that, our attack is identical over the polynomials in the vector of polynomials \mathbf{s} , as well as the polynomials in \mathbf{s}_1 and \mathbf{s}_2 and it is simply repeated to retrieve all polynomials in those secret vectors. Similarly, c^ℓ denotes the known polynomial, and \mathbf{c}^ℓ denotes its NTT representation.

It corresponds to the challenge polynomial c , which is among the output of the signature (Algorithm 1) for Dilithium. On the other hand, the public polynomial is any element of the \mathbf{u} part of the ciphertext for Kyber, which is an input of the key decapsulation (Algorithm 2). Again, our methodology is identical for any polynomial in \mathbf{u} .

Recall that, coefficient modulus $q = 3329$ and the ring dimension $n = 256$ for Kyber permits incomplete NTT of maximum $\log(n=2)$ layers. On the other hand, the specified coefficient modulus of Dilithium does not require the NTT to be incomplete. However, incomplete NTT can be preferred by the implementers to enhance performance [16]. Specifically, the multiplications $\mathcal{C}\mathbf{s}_1$ and $\mathcal{C}\mathbf{s}_2$ are referred to as *small NTT* and performed using a *carrier prime*. Although carrier primes are usually denoted by q^ℓ , we will use q to unify our notation. The small NTT operates with the prime $q = 257$ for Dilithium2 and Dilithium5 while $q = 769$ is chosen for Dilithium3. The rationale behind the small NTT is that, coefficients of $\mathcal{C}\mathbf{s}_1$ and $\mathcal{C}\mathbf{s}_2$ does not exceed n in absolute value. Recall that coefficient range of the polynomials in \mathbf{s}_1 and \mathbf{s}_2 is $[-n; n]$ while c has exactly n coefficients equal to 1 and the rest of the coefficients are 0.

The target implementation is the pqm4 library [17], which is mostly based on the work [16]. In the mentioned implementation, $\mathbf{r}_{[i:j]}$ (see Algorithm 3) are written to the memory while the intermediate steps of the computation remain in the processor. We consider the store instructions for $\mathbf{r}_{[i:j]}$ as the PoI, assuming a memory operation leads to a power leakage with a greater signal-to-noise-ratio (SNR) compared to the register updates or combinational logic. We use the HW model for the hypothetical power consumption computation.

VII. PROPOSED SIDE-CHANNEL ATTACKS

In this section, we first show a straightforward baseline attack and then give the details of a more efficient zero-value filtering attack. Figure 1 presents a high-level overview of the side-channel attacks discussed in this section.

A. The Baseline Attack

As each output coefficient $\mathbf{r}_{[i:j]}$, chosen as the PoI, depends on both $\mathbf{s}_{[i:0]}^\ell$ and $\mathbf{s}_{[i:1]}^\ell$, the *baseline scheme* is formed as a conventional non-profiled attack using one of the distinguishers presented in Section III that predicts $f_{\mathbf{s}_{[i:0]}^\ell; \mathbf{s}_{[i:1]}^\ell} \mathcal{G}$, simultaneously [21]. Consequently, q^2 hypotheses are tested by the baseline scheme.

We would like to make a note regarding the number of hypotheses for the attack on Dilithium. When the NTT multiplication function (Algorithm 3) is called, the coefficients $\mathbf{s}_{[i:j]}^\ell$ can be larger than the modulus q , due to the so-called lazy reductions. In particular, they are in the range $[-7q; 7q+]$, [16], [22], [33]. Fortunately, it is sufficient to predict in the set of residues, \mathbb{Z}_q , due to two main factors [21], [22]: 1) A trivial mathematical fact is that the coefficients are indeed in \mathbb{Z}_q , which is sufficient to compute the inverse NTT to obtain the coefficients of the secret polynomial. 2) For the selected PoI function, the integers in $[-7q; 7q+]$, that are in the same congruence class modulo q , mostly result in the same

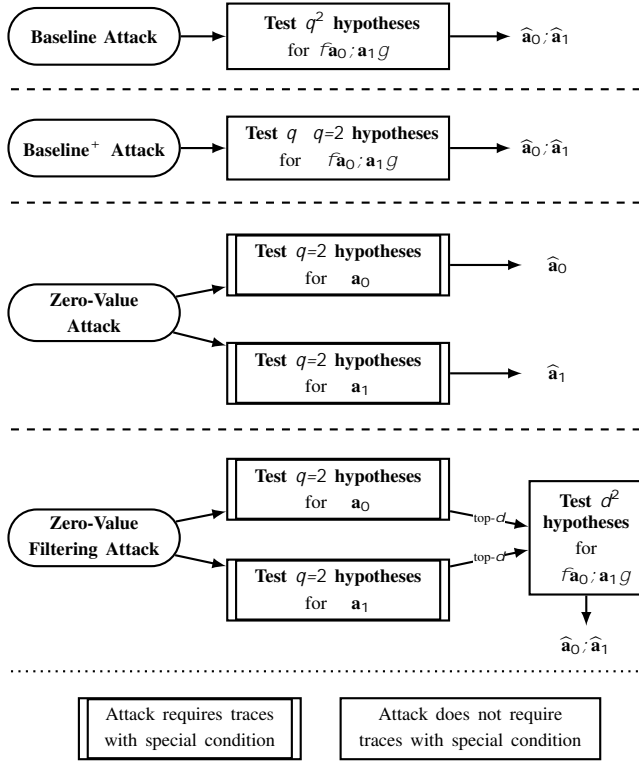


Fig. 1: Overview of the presented attacks. $f_{a_0}; a_1 g \in \mathbb{Z}_q^2$ denotes the attacked pair of secret coefficients, $f_{s_{[i,0]}}; s_{[i,1]} g$.

output. Therefore, the PoI that are calculated for integers of the same congruence class are correlated with each other. This situation does not exist in Kyber as the secret polynomials are stored in the NTT domain by algorithm definition so the coefficients are precisely in \mathbb{Z}_q .

The computational complexity of the baseline scheme is, therefore, $(q^2(n=2))$. q^2 is approximately 16:01-bit for Dilithium2 and Dilithium5, 19:17-bit for Dilithium3, while it is 23:4-bit for all security levels of Kyber. The baseline is an accurate, yet inefficient attacking scheme in terms of the attack run-time. Therefore, we seek more efficient methods to accelerate the attack time in the next section.

B. Using Negative Correlation

It is possible to further narrow down the hypothesis set presented in the preceding section by taking advantage of the negative correlation [22]. This is due to the fact that the Hamming weights of an integer and its additive inverse in 2 's complement notation are inversely correlated.

Note that, $\mathbf{V}_{:,f_{0;1}g} = \mathbf{V}_{:,f_{0;1}g}$, for any $f_{0;1}g \in \mathbb{Z}_q$; recalling that $\mathbf{V}_{:,f_{0;1}g}$ is the intermediate value vector computed w.r.t. $f_{0;1}g$. The statistical properties of Hamming Weight suggest that $\mathbf{H}_{:,f_{0;1}g}$ correlates with $\mathbf{H}_{:,f_{0;1}g}$. Therefore, $f_{s_{[i,0]}}; s_{[i,1]} g$ is retrieved by testing either the hypothesis set \mathbb{Z}_q or $\mathbb{Z}_{dq=2e}^+$, leading to q $q=2$ predictions.

Lastly, we need a distinguisher to tell the difference between the actual key and its additive inverse as the attacker can get either one of them. If the sign of the peak on correlation scores

is positive we conclude that the actual key is found. Otherwise, the additive inverse of the hypothesis is computed as the output of the attack. We would like to note that, the actual device leakage inversely correlates with the HW of intermediate data in some cases such as when the data-bus is pre-charged with all 1's. Then, the behavior of the explained distinguisher is reversed.

The improved baseline scheme that takes advantage of negative correlation is denoted by baseline⁺, dropping the attack complexity by 1-bit to $(q(q=2)(n=2))$.

C. Decreasing the Number of Hypotheses: Zero-Value Attack

A more practical scheme in terms of the attack run-time can be constructed by attacking the coefficients $s_{[i,0]}^0$ and $s_{[i,1]}^0$ individually, referred here as *Zero-Value (ZV) Attack*. To achieve this, we need to eliminate the effect of one of the secret coefficients from the reduction step during the NTT multiplication, whose output constitutes the chosen PoI. This can be accomplished by including only the traces to the attack that contain zeros in their coefficients, which multiply one of the secret coefficients. Consider line 2 of Algorithm 3 to develop intuition for the proposed method. Assume $c_{[i,1]}^0 = 0 \pmod q$ for some $0 < i < n-2$, then $s_{[i,1]}^0 c_{[i,1]}^0 \pmod q$ becomes 0 and $r_{[i,0]} = s_{[i,0]}^0 c_{[i,0]}^0 \pmod q$. With sufficient number of traces meeting the condition $c_{[i,1]}^0 = 0 \pmod q$, predictions on $s_{[i,0]}^0$ can be made independently from $s_{[i,1]}^0$ for the specific value of i . The prerequisites for the ZV attack are referred to as *zero-value conditions*. Table IV lists the four attacking scenarios that can be adopted. For instance, to attack $s_{[i,0]}^0$, we need the condition $c_{[i,1]}^0 = 0 \pmod q$ and use $c_{[i,0]}^0$ or $c_{[i,0]}^0 = 0 \pmod q$ and use $c_{[i,1]}^0$. As the conditions are identical with attack scenarios 1 and 4, both $s_{[i,0]}^0$ and $s_{[i,1]}^0$ will show peaks in the results.

Attacking Scenario	Target	Condition	Used Meta	Probability
1	$s_{[i,0]}^0$	$c_{[i,1]}^0 = 0$	$c_{[i,0]}^0$	0.0039
2	$s_{[i,0]}^0$	$c_{[i,0]}^0 = 0$	$c_{[i,1]}^0$	0.0039
3	$s_{[i,1]}^0$	$c_{[i,0]}^0 = 0$	$\delta_i c_{[i,1]}^0$	0.0039
4	$s_{[i,1]}^0$	$c_{[i,1]}^0 = 0$	$c_{[i,0]}^0$	0.0039

TABLE IV: ZV-Attack Scenarios. Probabilities are equal to $1/n$.

This approach leads to an attack complexity of (qn) and $((q=2)n)$ without and with the negative correlation trick, respectively, from the previous section. The drawback of this method is the hardness of finding traces meeting the mentioned conditions. We mark a trace and the corresponding known polynomial c as valid for the attack if at least one coefficient in c is 0; namely,

$$\begin{aligned} c_{[0,0]}^0 = 0 \text{ or } c_{[0,1]}^0 = 0 \text{ or } \dots \\ \dots c_{[n-2,0]}^0 = 0 \text{ or } c_{[n-2,1]}^0 = 0 \pmod q \end{aligned} \quad (7)$$

The probability for a random c^0 to be valid depends on q and n , and can be computed by the following

$$1 - \frac{q-1}{q}^n; \quad (8)$$

which leads to 0.283 for Dilithium and 0.074 for Kyber. Recall that, Dilithium's signature function applies rejection sampling, which means the target operation $\mathfrak{c}s_1$ is performed several times for each signature generation, with challenge polynomials that are thrown away since the corresponding signature is rejected. However, one can retrieve the unused challenge polynomials through another side-channel attack and include them in the attack to $\mathfrak{c}s_1$. We would like to note that, c is usually unprotected in the existing works from literature [34] [35] and the same assumption is made by [36]. As the NTT of c is computed, it can be attacked as presented in [20], which makes use of the belief-propagation algorithm to attack NTT transformation. Moreover, the coefficients of c are in $\mathbb{F}_{1;0;1}g$ so it would be relatively easier to distinguish between those. We leave this application as a future work. By multiplying the expected number of iterations presented in Table II by the probability 0.283, we find out that each signature operation contains 1.2, 1.44, 1.09 valid challenges, *i.e.* known polynomials, on average. On the other hand, \mathbf{u} is produced by the key encapsulation function of Kyber, whose inputs are publicly known. Therefore, an attacker can brute-force the seeds used by key encapsulation to find \mathbf{u} that satisfies the validity condition given in Equation 7.

The individual probabilities for the coefficients $c_{[i;j]}^0$ being 0 mod q , for a valid c^0 are another crucial factor of attack performance. Table IV lists the probabilities for the aforementioned conditions, which is $1/n = 0.0039$ for any i and j . Note that, each $\mathbf{s}_{[i;j]}^0$ is attacked with the ones ensuring the corresponding zero-value conditions among the collected traces. The listed probabilities suggest that the conditions are not met very often. Intuitively, assuming the SNR in \mathbf{T} requires 200 traces for the attack to converge, then the attacker must perform approximately 51.2K measurements on the victim's device considering the probabilities of conditions in Table IV. Although the attacking phase of the presented scheme is significantly faster than the baseline by a factor of $q=2$, a more optimal strategy exists, in terms of both the number of traces and attack run-time, as presented in the next section.

D. Decreasing the Number of Traces: Zero-Value Filtering

While the ZV scheme introduced in Section VII-C requires a large number of traces to retrieve the correct key exactly, alternatively having the correct key fall in top- d candidate list is relatively inexpensive in terms of the number of traces, depending on the value of d . Therefore, the ZV attack method can be used as a filtering mechanism for the following hypothesis testing, forming a two-stage attacking scheme, referred to as *Zero-Value Filtering Attack (ZV-FA)*, which is formalized in Figure 2.

In the first stage (a.k.a. *filtering stage*), $\mathbf{s}_{[i;0]}^0$ and $\mathbf{s}_{[i;1]}^0$ are attacked individually using the ZV attack scheme as presented in the preceding section. The outcomes of these ZV attacks are denoted by K^0 and K^1 , the sorted set of predictions for $\mathbf{s}_{[i;0]}^0$ and $\mathbf{s}_{[i;1]}^0$, respectively, based on the scores assigned by the employed distinguisher in ZV attacks. Then, in the second stage (a.k.a. *scoring stage*), a set of predictions $K = K_{[:d]}^0$ $K_{[:d]}^1$ of size d^2 is formed for the pair $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$. Notice

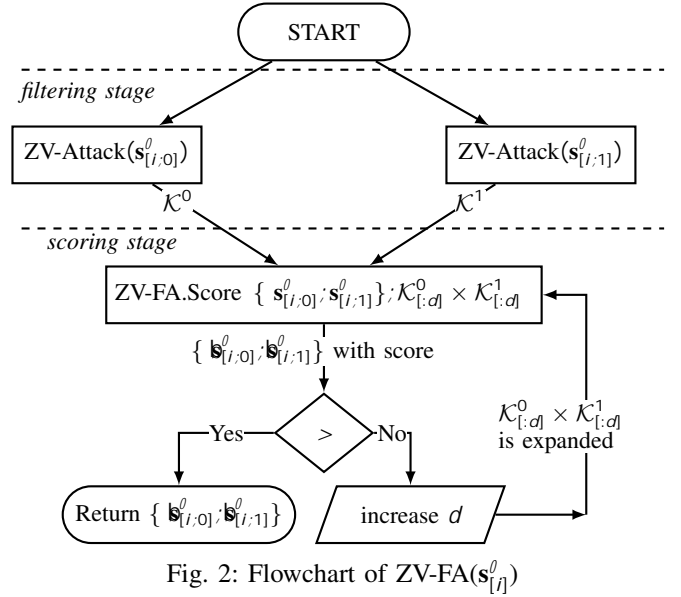


Fig. 2: Flowchart of ZV-FA($\mathbf{s}_{[i]}^0$)

that $K_{[:d]}^0$ ($K_{[:d]}^1$) stands for the top scoring d predictions in K^0 (K^1). Afterward, K is scored by one of the distinguishers presented in Section III. Compared to the baseline scheme, a relatively small number of hypotheses, d^2 is used for attacking $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$, as opposed to q^2 .

By the filtering stage, this method assumes that the correct key is in the top- d list of predictions of the highest scores for the ZV attack. A threshold mechanism, denoted by $>$, validates the assumption through the attack output. The value of d is iteratively increased and naturally $K_{[:d]}^0$ and $K_{[:d]}^1$ get larger, until a prediction scoring greater than $>$ is found. By increasing d , more candidates are evaluated by the second stage, which increases the probability of having the correct key in the top- d list; naturally, the second stage takes longer to evaluate more candidates. A trivial strategy for increasing d can be doubling it. Intuitively, doubling d is acceptable for Dilithium as q is relatively small, regarding the RAM usage and response times from scoring each set of candidates. However, reducing the rate of increasing d for Kyber can be desirable, as the search space can grow quite large, considering q^2 is 23.4-bit. We explicitly state how d is updated in our experiments in Section VIII. Needless to say, the evaluated candidates from previous trials are not included during the attack. The attack becomes identical to the baseline attack, for which the threshold is not taken into account if the scores remain below $>$ until d covers the whole search space.

The negative correlation trick presented in Section VII-B is also applied to ZV-FA both in the filtering stage and the scoring stage. Therefore, the number of hypotheses to be tested by ZV attacks in the filtering stage is $q=2$. Then, for each $\geq K^0$, we insert $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$ to K^0 , with the same rank as $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$. In this manner, either $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$ or its additive inverse $\mathfrak{f}\mathbf{s}_{[i;0]}^0; \mathbf{s}_{[i;1]}^0g$ is retrieved through the hypothesis testing in the scoring stage. The sign is corrected by observing the sign of the peak in correlation results as in the baseline⁺ scheme.

Compared to the ZV attack scheme, the new ZV-filtering attack is more effective with a significantly smaller number

of traces. The number of traces included in the filtering stage is denoted by N_F , while the second stage can be carried out without the zero-value conditions. As a result, it can be carried out with a sufficient number of traces to ensure that its output is reliable rather than using the entire set of valid traces, which is excessive for evaluating the score.

E. Improving ZV-FA by Using Inverse NTT to Validate Predictions

The zero-value filtering attack introduced in Section VII-D relies on the assumption that a precise threshold can be found for all attacks on $s_{[i]}$ for $0 < i < n=2$, which, however, may not hold in practice as a non-profiled attack is performed blindly. A possible solution to this problem is to use a conservative threshold. However, this approach is computationally expensive, and a conservative threshold can still result in false positives, albeit with a lower probability. Therefore, the attacker needs to verify the found secrets, s_1, s_2 for Dilithium, s for Kyber, by the MLWE equation presented in Section IV-A and Section IV-B, respectively. Note that this verification can only be performed after all the mentioned secrets have been attacked in all vector indices.

A more reasonable strategy for the attacker is to make use of the fact that $\text{NTT}^{-1}(s^\theta)$ is a short polynomial. Recall that, for Dilithium, $s_1 \in S^l, s_2 \in S^k$, whose coefficients are in the range $[-\beta; \beta]$. Similarly, for Kyber, s is sampled from B , for which the coefficients are in the range $[-\beta; \beta]$, as well. Therefore, a small mistake in the prediction will diffuse through the inverse NTT computation and ruin the coefficients of the output polynomial, empowering the attacker to efficiently validate the attack output.

Figure 3 illustrates the flowchart of the ZV-FA from a higher-level perspective with validation using the inverse NTT. Let \mathbf{s}^θ denote the vector of polynomials, a prediction to \mathbf{s}^θ , formed after completing the individual ZV attacks to $s_{[i;0]}^\theta$ and $s_{[i;1]}^\theta$ for all i at Step 1. To validate \mathbf{s}^θ , $\text{NTT}^{-1}(\mathbf{s}^\theta)$ is computed and the shortness property is sought in the resulting polynomial. If the found polynomial is not validated the mispredicted pair of coefficients in \mathbf{s}^θ is approximated and re-attacked to correct it. The approximation is performed by selecting the pair of coefficients with the minimum score. The current score for the prediction $\mathbf{s}_{[i]}^\theta$ is denoted by i . Observe that the outputs of ZV attacks are immediately scored at Step 2 if the shortness check fails. This initial scoring step is needed to be able to compare the ZV scores with ZV-FA scores, which are originally in distinct scales. The same distinguisher with ZV-FA.Score is applied to $\mathbf{s}_{[i]}^\theta$ using the top scorer candidates from Step 1, $K_{[0]}^{i;0}$ and $K_{[0]}^{i;1}$, with the same number of traces to get a comparable score with ZV-FA. In this manner, the attack terminates without re-attacking predictions which are already correctly predicted by ZV. Note that $K^{i;0}$ ($K^{i;1}$) is the set of predictions for $s_{[i;0]}^\theta$ ($s_{[i;1]}^\theta$) sorted for ZV attack scores, by slightly modifying our notation from the previous section as the coefficient index i is added to superscript.

The index of the minimum scoring pair from \mathbf{s}^θ is found by computing $i^\theta = \text{argmin}_i(i)$ and ZV-FA.Score is performed on $\mathbf{s}_{[i^\theta]}^\theta$ at Step 3 to replace $\mathbf{s}_{[i^\theta]}^\theta$. Here, we skip the filtering

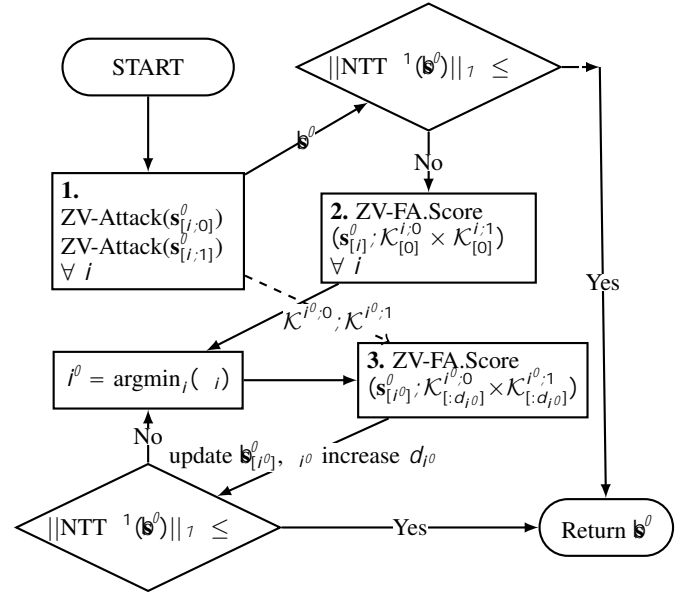


Fig. 3: ZV-FA for the whole vector of polynomials \mathbf{s}^θ with the application of inverse NTT validation

stage of ZV-FA (see Figure 2) because it is indeed performed at Step 1 and therefore $K^{i;0}$ and $K^{i;1}$ are known. $\mathbf{s}_{[i^\theta]}^\theta$ is updated if a better scoring prediction is found at the current invocation of ZV-FA.Score. At the same time, d_{i^θ} is updated as in Figure 2 (Again, the coefficient index i is included in the notation to differentiate between d for $\mathbf{s}_{[i]}^\theta$). Notice that the ZV-FA.Score can be performed for the same $\mathbf{s}_{[i]}^\theta$ multiple times. However, the scoring is performed with distinct d_i at each invocation of ZV-FA.Score to expand the search for the actual secret. Recall that from the previous section, predictions for $\mathbf{s}_{[i]}^\theta$ that are evaluated previously are not included in the attack. The prediction $\mathbf{s}_{[i]}^\theta$ is guaranteed to be corrected by subsequent applications of ZV-FA.Score if the correct value for $\mathbf{s}_{[i]}^\theta$ is the top-scoring among q ($q=2$) candidates.

The scheme is equivalent to ZV attack in terms of run-time if the found polynomial is validated after Step 1. On the other hand, the scheme evaluates $O(q(q=2)(n=2))$ predictions in the worst case, via iterations of Step 3, thus it becoming equivalent to the baseline⁺. The differentiation between both directions depends on the number of filtering traces, N_F . Note that the threshold presented in the previous section is not needed in the improved scheme thanks to inverse-NTT validation. The application of inverse NTT as a reliable and efficient method of verification renders the ZV-FA fault-tolerant. This method ensures the preservation of accuracy regardless of the choice of N_F , enhancing the attack performance.

VIII. RESULTS

In this section, we present the results obtained after implementing the above-mentioned attacks in a realistic experimental setting, against Dilithium-3 and Kyber-768. Moreover, we apply our attack against a masked version of Kyber-768.

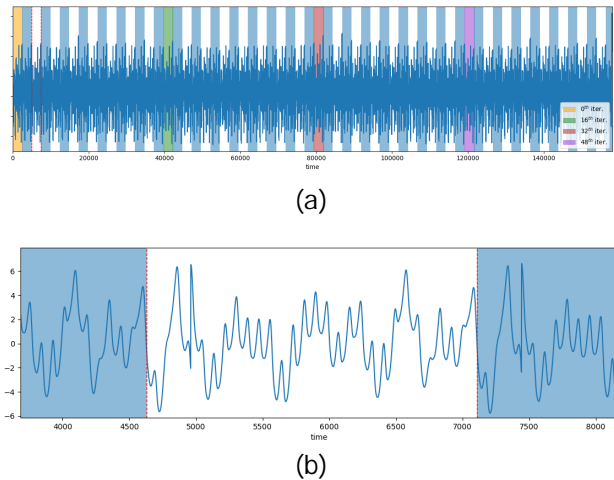


Fig. 4: Iterations of `_asymmetri c_mul _16_l oop`, implementation of Algorithm 3 for Dilithium, highlighted over mean trace

A. Experimental Setup

We employ Analog Devices' MAX32520¹ as the victim device to run pqm4's Dilithium signature and Kyber key decapsulation implementations. The first-order masked implementation of Kyber is developed on top of the pqm4's implementation by randomly splitting s^0 . The MAX32520 incorporates a 120 MHz ARM Cortex-M4F core that can sign random 32 B messages in 65.52 ms, on average, while it can perform Kyber's decapsulation in 7.2 ms. For EM trace collection, LeCroy WavePro HD oscilloscope² and Langer ICR HH500-6³ nearfield micro-probe are used. Sampling rate of the oscilloscope is set to 10 GS/s and 1 GS/s, yielding 83.33 and 8.33 samples per clock, for unprotected and protected implementations, respectively. We set up a trigger at the beginning of the implementation of Algorithm 3 from the target implementations for both Dilithium and Kyber to record the relevant time samples because we focus on the polynomial multiplication. The Scared library⁴ is used for analysis and attack, running on a computer equipped with 64GB RAM and AMD Ryzen 9 5900X 12-Core Processor clocked at 3.70 GHz.

B. Pre-processing and Analysis

To cope with the adverse effects of misalignment over time, we performed the following pre-processing steps for both algorithms: 1) pattern detection, 2) signal filtering, and 3) extraction around peaks. A reference pattern is set by band-pass filtering the first trace between 100 MHz and 140 MHz and applying moving variance to it. The traces are aligned based on the reference pattern. Then, 64 peaks (128 for first-order protected Kyber), which correspond to iterations of Algorithm 3 (loop is unrolled by a factor of two), are detected and sequential points after each peak are

combined. Figure 4 highlights the iterations over the average of pre-processed traces for Dilithium, conforming with the pre-knowledge on the implementation. On the other hand, iterations of Algorithm 3 for (masked) Kyber, are observed (for both shares) in Figure 5. Given the clear visibility of the iterations of Algorithm 3 over time samples, it is possible to conduct individual attacks on $s_{[i]}^0$ in time regions associated with each iteration. Note that, partitioning the attack range over time is critical for the presented performance results of all schemes.

Upon observation, we use the concatenation of $\mathbf{r}_{[i,0]}$ and $\mathbf{r}_{[i,1]}$ (see Algorithm 3) as the PoI for attacking Dilithium, while we use $\mathbf{r}_{[i,0]}$ for Kyber. As an initial analysis, we performed the `baselineCPA` on $s_{[0]}^0$ and $s_{[1]}^0$. The convergence patterns of the retrieved secrets are presented in Figure 6.

C. Attack and Performance

1) *First-order*: For the first-order attacks that target unprotected implementations of Dilithium and Kyber, we use CPA as the distinguisher. We start the evaluation by the performance of the `baselineCPA` and `baselineCPA+` schemes. Figure 7 illustrates the distribution of the required number of traces $s_{[i]}^0$ needed to converge in our experiments. Note that, the histograms are computed based on a single s^0 and varying i . Accordingly, 220 and 150 traces are needed the retrieve whole s^0 for Dilithium and Kyber, respectively. Notice that, these numbers are the maximums of the values presented in the histograms. However, a significant portion of the secret coefficients, $s_{[i]}^0$, are indeed revealed with fewer traces, around 50. The performance of the `baselineCPA` and `baselineCPA+` schemes is reported in Table V. In terms of accuracy, both schemes exhibit flawless performance and do not pose any concerns regarding accuracy. However, in terms of run-time, the performance of the attacks is moderate. As expected, the `baseline+` scheme improves the performance of the baseline scheme by a factor of 2 while preserving accuracy, which supports the correctness of the attack methodology. Nevertheless, even with the `baseline+` scheme, retrieving s_1 and s_2 requires approximately 4.37 hours for Dilithium while it takes 22.4 hours to attack s for the Kyber case. Note that the `baselineCPA` is equivalent to the attack presented by [21] against the same implementation of Kyber. Our application of the conventional approach is slightly better than that work both in terms of attack run-time and number of traces.

For the application of ZV-FA to Dilithium (to Kyber), the attack scenarios 1 and 2 (scenario 1 for Kyber) from Table IV are employed for attacking the lower degree coefficients of the polynomials, specifically $s_{[i,0]}^0$ for any $0 \leq i < 128$, while the scenarios 3 and 4 (scenario 3) are utilized for the higher degree coefficients $s_{[i,1]}^0$. It should be noted that scenarios 1 and 4 are not independently executed, as they represent the same attack. The outcomes of different scenarios are combined by multiplying their respective results. For both schemes, we use $N = 500$ traces for ZV-FA.Score, and d_i is doubled to increase it after each call to ZV-FA.Score (see Figure 3). Note that, we use slightly more traces compared to the convergence of the baselines (see Figure 6). This is needed to have the score

¹<https://www.analog.com/en/products/max32520.html>

²<https://teledynelecroy.com/oscilloscope/wavepro-hd-oscilloscope>

³<https://www.langer-emv.de/en/product/near-field-microprobes-icr-hh-h-field/26/icr-hh500-6-near-field-microprobe-2-mhz-to-6-ghz/108>

⁴<https://pypi.org/project/scared/>

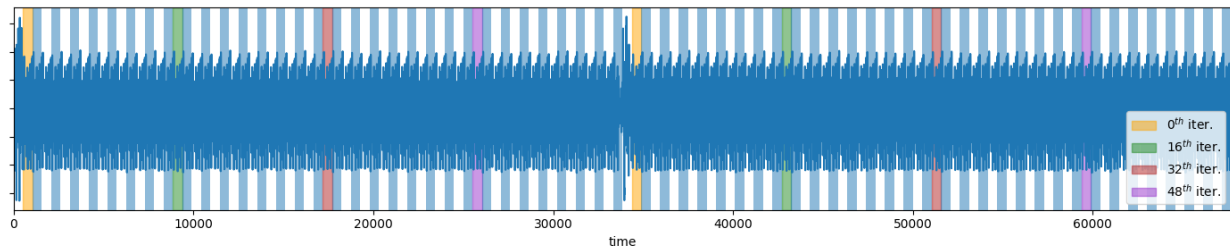


Fig. 5: Iterations of `frombytes_mul_asm_acc_32_16` for two shares, implementation of Algorithm 3 for Kyber, highlighted over mean trace.

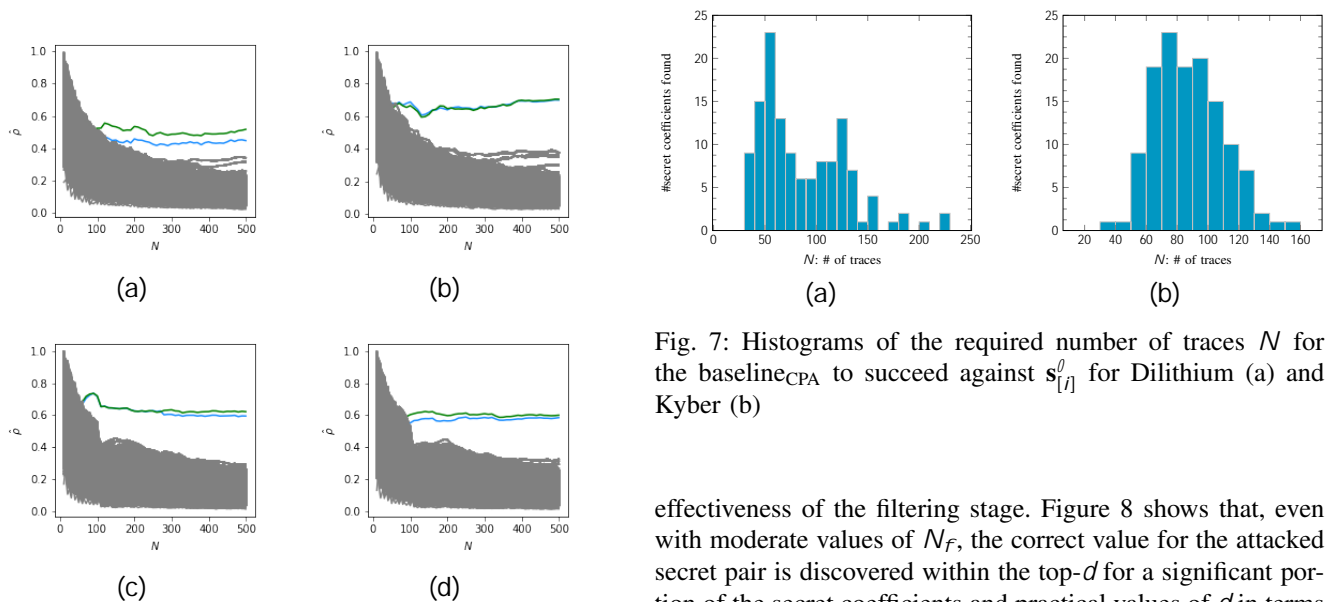


Fig. 6: Key convergence of $\text{Baseline}_{\text{CPA}}$, for $s_{[0]}^{\ell}$ (a,c) $s_{[1]}^{\ell}$ (b,d) for Dilithium (a,b) and Kyber (c,d). Green lines denote the correct hypothesis while the blue line denotes its additive inverse

of $s_{[i]}^{\ell}$ converged and become comparable with the scores of predictions to other secret coefficients. Observe from Figure 6 that the SNR of even and odd coefficients are different in Dilithium. Therefore we scale the scores onto the same range by multiplying the odd coefficients by $\alpha = 1/5 = 0.2$.

Thanks to the inverse NTT validation and correction mechanism presented in Section VII-E, The ZVFA's accuracy is independent of the number of traces used for filtering $N_{\mathcal{F}}$, which, however, determines its performance. Recall that, the performance of the ZV-FA scheme strongly depends on the

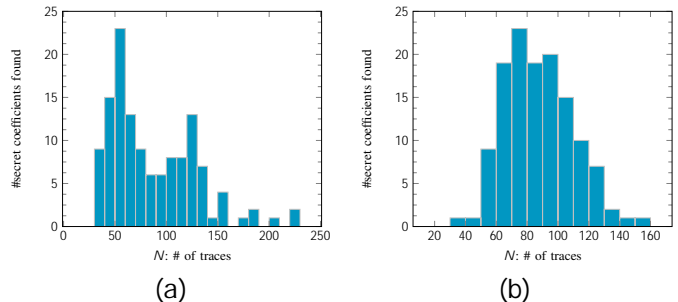


Fig. 7: Histograms of the required number of traces N for the $\text{baseline}_{\text{CPA}}$ to succeed against $s_{[i]}^{\ell}$ for Dilithium (a) and Kyber (b)

effectiveness of the filtering stage. Figure 8 shows that, even with moderate values of $N_{\mathcal{F}}$, the correct value for the attacked secret pair is discovered within the top- d for a significant portion of the secret coefficients and practical values of d in terms of performance. Particularly for Dilithium and $N_{\mathcal{F}} = 5\text{K}$, 203 of 256 (%80) secret coefficients are retrieved in top-64, which corresponds to $(769 - 64) = 705$ (%92) reduction in the search space from the baseline to ZV-FA.Score. Similarly for Kyber and the same value of $N_{\mathcal{F}}$, 202 coefficients are in the top-256, leading to (%92) reduction in the search space.

Experimental results indicate that the ZV filtering can substantially decrease the attack response time up to three orders of magnitude, depending on the number of filtering traces $N_{\mathcal{F}}$ available in the system. The trade-off between $N_{\mathcal{F}}$ and speed-up is illustrated in Figure 9 for both attacked algorithms. Observe that the trade-off suggests the same pattern for Dilithium and Kyber. The ZV-FA approaches to the ZV attack as $N_{\mathcal{F}}$ increases and approaches to the baseline as $N_{\mathcal{F}}$ decreases. Notably, even a small number of traces can significantly improve baseline performance. For example with

Algorithm	Method	N	Runtime($s_{[i]}^{\ell}$)	Runtime(s^{ℓ})	
					Runtime(s_1, s_2)
Dilithium-3	$\text{Baseline}_{\text{CPA}}$	220	22.35s	48m	8.74h
Dilithium-3	$\text{Baseline}_{\text{CPA}}^+$	220	11.18s	24m	4.37h
					Runtime(s)
Kyber-768	$\text{Baseline}_{\text{CPA}}$	150	7m	14.9h	44.8h
Kyber-768	$\text{Baseline}_{\text{CPA}}^+$	150	3.5m	7.45h	22.4h
Kyber	$\text{Baseline}_{\text{CPA}}$ [21]	200	5m	10.7h	

TABLE V: Performance of Baseline and Baseline⁺ Attacks

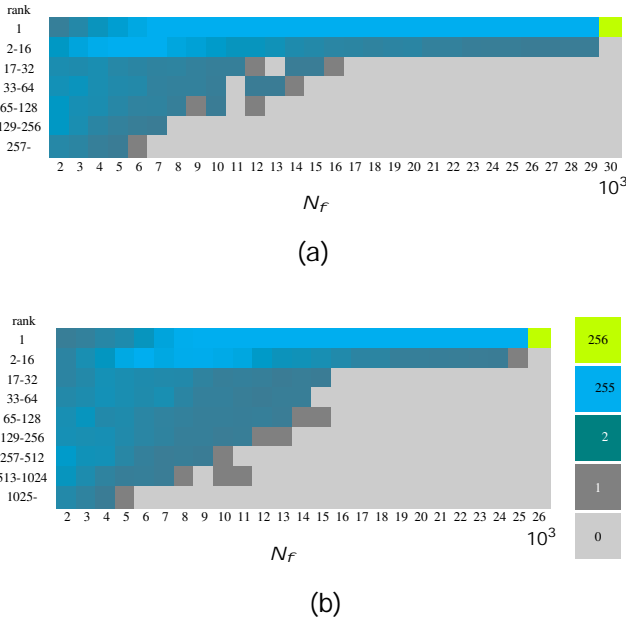


Fig. 8: Histograms of the ranks of correct hypotheses for $s_{[i;j]}^0$ in $K^{i;j}$ during filtering stage of ZV-FA, *i.e.* the correct secret among results of ZV attacks, w.r.t. N_F for Dilithium (a) and Kyber (b)

$N_F = 5K$ the collection of which is feasible, the ZV-FA provides a speed-up of 17 and 30 for Dilithium and Kyber, respectively.

When more valid traces are available in the system, particularly with $N_F = 18K$, ZV-FA achieves a speed-up of 362 for Dilithium over the $\text{baseline}_{\text{CPA}}$. We underline that, the achieved speed-up can save approximately 523 minutes (8:71 hours) considering the retrieval of whole s_1 and s_2 . Recall that $k = 6$ and $l = 5$ for Dilithium-3, which means s_1 and s_2 consist of 5 and 6 secret polynomials, respectively. As for Kyber (recall that $k = 3$ for Kyber-768, which means s has 3 elements.), our scheme is more favorable as q is roughly 2-bit larger compared to Dilithium-3. With $N_F = 17K$, ZV-FA achieves **1915** speed-up over the $\text{baseline}_{\text{CPA}}$. It saves roughly 44.77 hours of computation time, considering all the elements of Kyber’s secret vector of polynomials s .

On the other hand, ZV-FA reduces the number of traces needed for the ZV attack while the run-time performance is slightly improved. Observe that the ZV attack is successful with $N_F = 30K$ for Dilithium which brings up a speed-up of 313. The same speed-up is achieved by ZV-FA with $N_F = 14K$. Similarly for Kyber, the ZV attack is successful with $N_F = 26K$ accelerating the $\text{baseline}_{\text{CPA}}$ by 1508 while ZV-FA reaches the same speed-up with $N_F = 14.5K$.

Recall also that another study [22] targets Dilithium with a non-profiled attack. However, the target implementation uses the original 23-bit coefficient modulus of Dilithium. Therefore, our study is not comparable to [22] as the target implementations are different. On the other hand, while their method accelerates the baseline approach about 16 times, ours provides a speedup of more than two orders of magnitude.

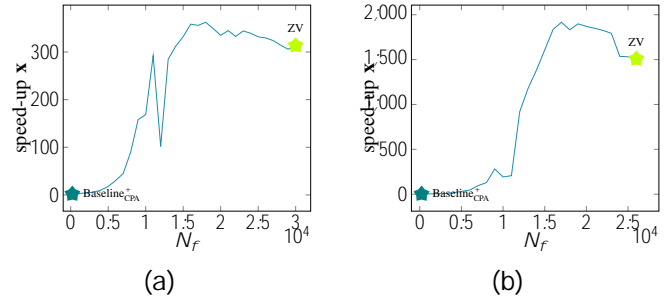


Fig. 9: Speed-up of ZV-FA w.r.t. N_F for Dilithium (a) and Kyber (b). The performance of $\text{Baseline}_{\text{CPA}}^+$ and ZV attack are marked.

Method	i	N	Runtime($s_{[i]}^0$)	Runtime($s_{[i]}^0$) $n/2$ k Runtime(s)
$\text{Baseline}_{\text{HOCPA}}^+$	0	14K	78m	21d
$\text{Baseline}_{\text{HOCPA}}^+$	1	23K	129m	34d
$\text{Baseline}_{\text{MMIA}}^+$	0	7K	182m	48.5d

TABLE VI: Performance of Baseline and Baseline⁺ Attacks against first-order protected Kyber-768

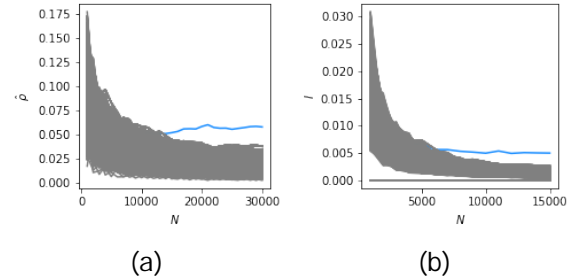


Fig. 10: Key convergence for $s_{[0]}^0$ $\text{Baseline}_{\text{HOCPA}}^+$ (a) $\text{Baseline}_{\text{MMIA}}^+$ (b) for masked Kyber. Blue line denotes the additive inverse of the correct hypothesis

Consequently, we expect their method would not be as effective as ours when an incomplete NTT method is used in the implementation.

2) *Second-order*: To discuss the efficiency of our attack in the protected case, we first present the performance of the baseline schemes thereof in Table VI. Due to long running times, we perform the evaluation only for $s_{[0]}^0$ and $s_{[1]}^0$. The last column approximates the required amount of time to break whole secret key s , based on the statistics of the attacks to $s_{[0]}^0$ and $s_{[1]}^0$. For instance, retrieving s would roughly take 68 days if all coefficients were retrieved by 23K traces as $s_{[1]}^0$. Therefore the baseline scheme stands as an impractical option. The speed-up of ZV-FA is computed based on the average number of traces needed by $\text{Baseline}_{\text{HOCPA}}$ over the analyzed coefficients, $(23 + 14) \cdot 2 = 18.5K$.

Observe from Table VI and Figure 10 that MMIA is superior to HOCPA in terms of the number of traces while HOCPA runs faster. Therefore, we use MMIA as the distinguisher in the filtering stage, differently from the unprotected case. We

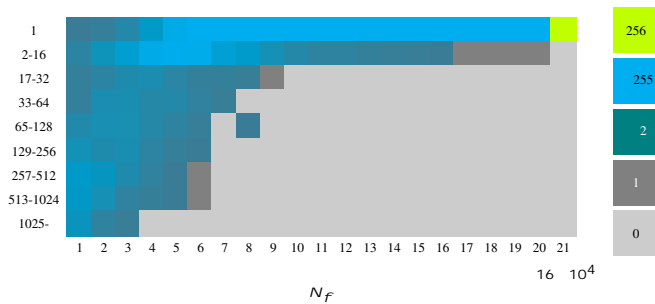


Fig. 11: Histograms of the ranks of the correct hypotheses for $s_{[i,j]}^0$ in $K^{i,j}$ during filtering stage of ZV-FA, w.r.t. N_F , against masked Kyber

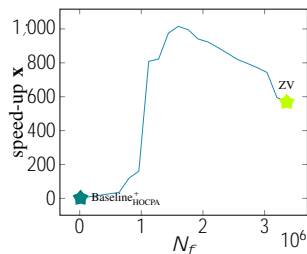


Fig. 12: Speed-up of ZV-FA w.r.t. N_F against Masked Kyber.

utilize MMIA with two bins and select the bins such that it partitions the samples into equal parts, *i.e.* halves. As the iterations of Algorithm 3 are easily distinguishable through the mean trace as depicted by Figure 5, we use constant offset for combining points over time samples. We would like to note that we observe that the MMIA is quite efficient in this setting. However, we leave a comprehensive study on MMIA in comparison with HOCPA regarding arithmetic masking as a future work. On the other hand, we use HOCPA in the scoring stage, considering its superior run-time performance and we employ $N = 50K$ traces for scoring. Recall that, significantly more traces are available during scoring as zero-value conditions are not sought therein.

Figure 11 demonstrates the distribution of the ranking of the correct hypothesis for $s_{[i,j]}^0$, among the sorted results from the filtering stage of ZV-FA, $K^{i,j}$. We observe that the filtering stage is quite effective as in the unprotected case. Particularly with $N_F = 480K$, for 213 out of 256 secret coefficients, the correct hypothesis is in top 256, *i.e.* $K_{[1:256]}^{i,j}$, leading to %92 reduction in search space. Figure 12 depicts the speed-up values achieved over the Baseline_{HOCPA}. We observe that ZV-FA reaches **1015** speed-up over Baseline_{HOCPA}, with $N_F = 1.6M$. If those many traces are not available to the attacker, ZV-FA is still more practical compared to the Baseline_{HOCPA}. For instance, a speed-up of 35 is observed with a relatively less number of filtering traces, $N_F = 640K$.

IX. CONCLUSION AND FUTURE WORK

This paper presents a series of non-profiled side-channel attacks against the incomplete NTT-based polynomial multiplication (Algorithm 3), which is widely adopted in lattice-based cryptography. We exercised our approach in the proposed

attacks against the signature algorithm Dilithium and the KEM Kyber [16], [17]. Specifically, the attacks focus on the NTT-based polynomial multiplications $\mathfrak{C}\mathfrak{s}_1$ and $\mathfrak{s}^T\mathfrak{u}$. The target implementation for Dilithium operates with a carrier prime $q^0 = 769$, which restricts the NTT to be incomplete. On the other hand, Kyber also needs to employ 7 layers of incomplete NTT by algorithm definition, albeit with a different prime $q = 3329$. The baseline and baseline+ schemes are conventional methods that rely on brute-force methods in the sets of cardinality q^2 and $q^2=2$, respectively, as two coefficients of the incomplete NTT representation must be predicted together.

To mitigate the search costs, we introduced the zero-value attack, which reduces the size of the set of hypotheses to q in the brute force attack by taking advantage of multiplication by 0 to eliminate one of the attacked pair of coefficients from the equation. However, this approach requires a significantly higher number of traces. Next, we presented the zero-value filtering attack, which represents a trade-off between the number of traces and attack run-time. With an appropriate number of traces, this attack can achieve a speed-up of two orders of magnitude over the baseline. Finally, we proposed an efficient way of verification of predictions on short polynomials, utilizing the inverse NTT transformation. It makes the proposed scheme accurate independent of the number of filtering traces. Experiments suggest that the ZV-FA is favorable even with moderate parameters. Moreover, we show that our attack is effective in the presence of masking by applying it against a first-order protected implementation of Kyber. Our approach is also generalizable to higher orders as long as the attacker can combine leaky samples over time that correspond to the different shares. During the study, we found out that MMIA outperforms HOCPA in terms of the number of traces and we leave this study as a future work.

ACKNOWLEDGMENTS

This work is partially supported by the European Union's Horizon Europe research and innovation program under grant agreement No: 101079319 and by TUBITAK under Grant Number 122E222.

REFERENCES

- [1] N. Koblitz, "Elliptic curve cryptosystems," *Mathematics of computation*, vol. 48, no. 177, pp. 203–209, 1987.
- [2] V. S. Miller, *Use of elliptic curves in cryptography*. Springer, 1986.
- [3] C.-P. Schnorr, "Efficient identification and signatures for smart cards," in *Advances in Cryptology—CRYPTO'89 Proceedings 9*. Springer, 1990, pp. 239–252.
- [4] R. L. Rivest, A. Shamir, and L. Adleman, "A method for obtaining digital signatures and public-key cryptosystems," *Communications of the ACM*, vol. 21, no. 2, pp. 120–126, 1978.
- [5] P. W. Shor, "Algorithms for quantum computation: discrete logarithms and factoring," in *Proceedings 35th annual symposium on foundations of computer science*. Ieee, 1994, pp. 124–134.
- [6] S. Bai, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, P. Schwabe, G. Seiler, and D. Stehlé, "Crystals-dilithium algorithm specifications and supporting documentation (version 3.1)," *NIST Post-Quantum Cryptography Standardization Round*, vol. 3, 2021.
- [7] O. Regev, "On lattices, learning with errors, random linear codes, and cryptography," *Journal of the ACM (JACM)*, vol. 56, no. 6, pp. 1–40, 2009.

