

A Side-Channel Attack on a Hardware Implementation of CRYSTALS-Kyber

Yanning Ji
KTH Royal Institute of Technology
Stockholm, Sweden
yanning@kth.se

Ruize Wang
KTH Royal Institute of Technology
Stockholm, Sweden
ruize@kth.se

Kalle Ngo
KTH Royal Institute of Technology
Stockholm, Sweden
kngo@kth.se

Elena Dubrova
KTH Royal Institute of Technology
Stockholm, Sweden
dubrova@kth.se

Linus Backlund
KTH Royal Institute of Technology
Stockholm, Sweden
lbackl@kth.se

Abstract—CRYSTALS-Kyber has been recently selected by the NIST as a new public-key encryption and key-establishment algorithm to be standardized. This makes it important to assess how well CRYSTALS-Kyber implementations withstand side-channel attacks. Software implementations of CRYSTALS-Kyber have been already analyzed and the discovered vulnerabilities were patched in the subsequently released versions. In this paper, we present a profiling side-channel attack on a hardware implementation of CRYSTALS-Kyber with the security parameter $k = 3$, Kyber768. Since hardware implementations carry out computation in parallel, they are typically more difficult to break than their software counterparts. We demonstrate a successful message (session key) recovery by deep learning-based power analysis. Our results indicate that currently available hardware implementations of CRYSTALS-Kyber need better protection against side-channel attacks.

Index Terms—Post-quantum cryptography, CRYSTALS-Kyber, LWE-based KEM, side-channel attack, FPGA, power analysis, deep learning

I. INTRODUCTION

In July 2022, the National Institute of Standards and Technology (NIST) has selected CRYSTALS-Kyber as a new public-key encryption and key-establishment algorithm to be standardized [1]. CRYSTALS-Kyber is a lattice-based cryptographic algorithm that is expected to be capable of protecting confidential information after the emergence of large-scale quantum computers. Even if it might take many years until large-scale quantum computers become a reality, the need for long-term security necessitates an early start of the transition. The industry is preparing to plan and budget for a shift to quantum-resistant cryptographic algorithms. The National Security Agency (NSA) has recently included CRYSTALS-Kyber in the suite of cryptographic algorithms recommended for national security systems [2].

The standardization of CRYSTALS-Kyber makes it important to assess the resistance of its implementations to side-channel attacks (SCAs). Several software implementations of CRYSTALS-Kyber have been already analyzed [3]–[6] and the discovered vulnerabilities were patched in the subsequently released versions. However, to the best of our

knowledge, no successful attack on a hardware implementation of CRYSTALS-Kyber has been reported until now. Hardware implementations are typically significantly more difficult to break than software variants because they carry out the computation in parallel. There are opinions that the parallelism, in a combination with a smaller minimum feature size of application-specific integrated circuits (ASICs) of field-programmable gate array (FPGAs) process technologies may provide a sufficient resistance against side-channel analysis without any SCA-specific countermeasures.

Our Contributions: In this paper, we present a profiling side-channel attack on a hardware implementation of CRYSTALS-Kyber with the security parameter $k = 3$, Kyber768. Using deep learning-based power analysis, we successfully recover messages from a Xilinx Artix-7 FPGA implementation of CRYSTALS-Kyber from [7]. In CRYSTALS-Kyber, a message recovery from a properly generated ciphertext implies the session key recovery since the session key is derived from the message using hash functions. Furthermore, by recovering messages contained in seven chosen ciphertexts [8], one can extract the secret key of CRYSTALS-Kyber. Unlike the non-profiling attack on Kyber512 presented in [9], our attack does not require any knowledge of register reference values.

The success of our attack is due to the new *sliced multi-bit error injection* method which is an extension of the original multi-bit error injection method introduced in [10] for side-channel analysis of software implementations of lattice-based PKE/KEMs. Software implementations do not require slicing because they execute instructions sequentially. However, slicing is essential for hardware implementations. Our experimental results show that, without slicing, only 10% of messages can be recovered with enumeration up to 2^{64} . With slicing, we can recover all messages with the same enumeration.

The rest of this paper is organized as follows. Section II reviews previous work related to side-channel analysis of CRYSTALS-Kyber implementations. Section III gives a background on CRYSTALS-Kyber algorithm. Section IV describes the equipment used in the experiments. Sections V and VI

present the profiling and attack stages, respectively. Section VII summarizes experimental results. Section VIII concludes the paper and discusses future work.

II. PREVIOUS WORK

Side-channel attacks on implementations of post-quantum cryptographic algorithms have drawn a lot of attention since the beginning of the NIST post-quantum cryptography (PQC) standardization process in 2016. The lattice-based public key encryption (PKE) algorithms and key encapsulation mechanisms (KEMs) were the most popular targets. This is due to the fact that many of the NIST PQC candidates are based on lattice problems: an NTRU-based scheme NTRU [11], a Learning With Errors (LWE)-based scheme CRYSTALS-Kyber [12], and a Learning With Rounding (LWR)-based scheme Saber [13]. The majority of the presented side-channel attacks are on software implementations.

In [14], profiling deep-learning-based message recovery attacks using a single power trace representing the execution of the encapsulation procedure are demonstrated for unprotected software implementations of CRYSTALS-Kyber, Saber, and FrodoKEM. In [15], near field EM secret key recovery attacks on unprotected and protected implementations of the same algorithms are presented. It is shown how their masked implementations can be broken in two steps, by attacking each share individually. In [16] an attack on a first-order masked implementation of CRYSTALS-Kyber using the one-step method introduced in [17] is reported. It exploits the message encoding vulnerability discovered in [14].

In [18], secret key recovery attacks on software implementations of lattice-based KEMs, including CRYSTALS-Kyber, Saber, and NTRU, based on power/near field EM analysis are demonstrated. These attacks make use of the side-channel leakage during the re-encryption step of the decapsulation procedure. In [19], correlation power analysis-based attacks of all lattice-based candidates of the NIST are presented, targeting the polynomial multiplication in unprotected software implementations.

In [8], a near field EM emanations-based message recovery attack on an unprotected software implementation of CRYSTALS-Kyber making use of a vulnerability in the Fujisaki-Okamoto (FO) transform is described. In [20] an EM-based chosen ciphertext side-channel attack on an unprotected software implementation of CRYSTALS-Kyber is shown. In [21], side-channel attacks on two implementations of masked polynomial comparison, demonstrated on CRYSTALS-Kyber, are presented.

All attacks mentioned above targeted software implementations of CRYSTALS-Kyber and other lattice-based algorithms. Since these algorithms perform many repeated computations, it is of advantage to use hardware for their implementations. In [22] a hardware implementation of CRYSTALS-Kyber in FPGA is described which takes 10 times less clock cycles than the ARM Cortex-M4 implementation of CRYSTALS-Kyber from [23].

In [7] an even smaller and faster hardware implementation of CRYSTALS-Kyber is presented. With suitably designed pipelines and well-optimized architecture, this implementation is capable of executing the decapsulation procedure of any version of CRYSTALS-Kyber in 14,000 clock cycles while consuming 7,500 LUTs only. Thus, it can fit into the smallest Xilinx Artix-7 series devices.

However, the evaluation of resistance of hardware implementations of PQC KEMs to side-channel analysis is still in the beginning. In [24] a power-based key recovery attack on the McEliece cryptosystem, a round 4 candidate of the NIST PQC, implemented in Xilinx Artix-7 FPGA is presented. In [25], a deep learning-based power analysis of Saber implemented in Xilinx Artix-7 FPGA is presented. Due to the parallel processing of hardware, the MLP models are able to predict the Hamming weight of messages with the accuracy of 88% at most. The message bits cannot be distinguished accurately. In [26] a preliminary power analysis of a hardware implementation of CRYSTALS-Kyber from [22] is presented. Leakage points were found using a t-test for the Kyber512 implementation in Xilinx Virtex-7 FPGA.

In [9], a non-profiling correlation EM analysis of an FPGA implementation of Kyber512 (security parameter $k = 2$) is presented targeting polynomial multiplication during the PKE decryption step of the decapsulation algorithm. The attack utilizes 166,620 traces to recover the secret key. However, it requires knowledge of the register reference values.

To the best of our knowledge, no successful side-channel attack on a hardware implementation of Kyber768 (security parameter $k = 3$) has been demonstrated so far. Typically, the higher the security parameter, the more difficult the implementation is to break.

III. CRYSTALS-KYBER ALGORITHM

CRYSTALS-Kyber [12] is an IND-CCA2-secure cryptographic algorithm which means that it is indistinguishable under an adaptive chosen ciphertext attack. The security of CRYSTALS-Kyber relies on the hardness of the Module Learning with Errors (M-LWE) problem.

CRYSTALS-Kyber contains a CPA-secure PKE scheme, KYBER.CPAPKE, and a CCA-secure KEM scheme, KYBER.CCAKEM, based on a post-quantum version of the Fujisaki-Okamoto transform [27]. These algorithms are described in Fig. 1 and Fig. 2 respectively.

Let \mathbb{Z}_q denote the ring of integers modulo a prime q , and R_q denote the quotient ring $\mathbb{Z}_q[X]/(X^n + 1)$. CRYSTALS-Kyber works with vectors of ring elements in R_q^k , where k is the lattice dimension used to scale the security levels.

The term $x \leftarrow \chi(S; r)$ denotes sampling x from a distribution χ over a set S using seed r . The uniform distribution is denoted by \mathcal{U} . The centered binomial distribution with parameter μ is denoted by β_μ .

The Decode_l function is designed to decode an array of $32l$ bytes into a 256 length polynomial with each coefficient in $\{0, 1, \dots, 2^l - 1\}$. The Encode_l function is the inverse of Decode_l , which encodes each polynomial coefficient

KYBER.CPAPKE.KeyGen()

- 1: $(\rho, \sigma) \leftarrow \mathcal{U}(\{0, 1\}^{256})$
- 2: $\mathbf{A} \leftarrow \mathcal{U}(R_q^{k \times k}; \rho)$
- 3: $\mathbf{s}, \mathbf{e} \leftarrow \beta_{\eta_1}(R_q^{k \times 1}; \sigma)$
- 4: $\mathbf{t} = \text{Encode}_{12}(\mathbf{A}\mathbf{s} + \mathbf{e})$
- 5: $\mathbf{s} = \text{Encode}_{12}(\mathbf{s})$
- 6: **return** $(pk = (\mathbf{t}, \rho), sk = \mathbf{s})$

KYBER.CPAPKE.Dec(\mathbf{s}, c)

- 1: $\mathbf{u} = \text{Decompress}_q(\text{Decode}_{d_u}(c_1), d_u)$
- 2: $v = \text{Decompress}_q(\text{Decode}_{d_v}(c_2), d_v)$
- 3: $\mathbf{s} = \text{Decode}_{12}(\mathbf{s})$
- 4: $m = \text{Encode}_1(\text{Compress}_q(v - \mathbf{s}^T \mathbf{u}, 1))$
- 5: **return** m

KYBER.CPAPKE.Enc($pk = (\mathbf{t}, \rho), m, r$)

- 1: $\mathbf{t} = \text{Decode}_{12}(\mathbf{t})$
- 2: $\mathbf{A} \leftarrow \mathcal{U}(R_q^{k \times k}; \rho)$
- 3: $\mathbf{r} \leftarrow \beta_{\eta_1}(R_q^{k \times 1}; r)$
- 4: $\mathbf{e}_1 \leftarrow \beta_{\eta_2}(R_q^{k \times 1}; r); e_2 \leftarrow \beta_{\eta_2}(R_q^{1 \times 1}; r)$
- 5: $\mathbf{u} = \mathbf{A}^T \mathbf{r} + \mathbf{e}_1$
- 6: $v = \mathbf{t}^T \mathbf{r} + e_2 + \text{Decompress}_q(m, 1)$
- 7: $c_1 = \text{Encode}_{d_u}(\text{Compress}_q(\mathbf{u}, d_u))$
- 8: $c_2 = \text{Encode}_{d_v}(\text{Compress}_q(v, d_v))$
- 9: **return** $c = (c_1, c_2)$

Fig. 1: Description of KYBER.CPAPKE algorithm from [12].

KYBER.CCAKEM.KeyGen()

- 1: $z \leftarrow \mathcal{U}(\{0, 1\}^{256})$
- 2: $(pk, \mathbf{s}) = \text{KYBER.CPAPKE.KeyGen}()$
- 3: $sk = (\mathbf{s}, pk, \mathcal{H}(pk), z)$
- 4: **return** (pk, sk)

KYBER.CCAKEM.Encaps(pk)

- 1: $m \leftarrow \mathcal{U}(\{0, 1\}^{256})$
- 2: $m = \mathcal{H}(m)$
- 3: $(\hat{K}, r) = \mathcal{G}(m, \mathcal{H}(pk))$
- 4: $c = \text{KYBER.CPAPKE.Enc}(pk, m, r)$
- 5: $K = \text{KDF}(\hat{K}, \mathcal{H}(c))$
- 6: **return** (c, K)

KYBER.CCAKEM.Decaps($sk = (\mathbf{s}, pk, \mathcal{H}(pk), z), c$)

- 1: $m' = \text{KYBER.CPAPKE.Dec}(\mathbf{s}, c)$
- 2: $(\hat{K}', r') = \mathcal{G}(m', \mathcal{H}(pk))$
- 3: $c' = \text{KYBER.CPAPKE.Enc}(pk, m', r')$
- 4: **if** $c = c'$ **then**
- 5: **return** $K = \text{KDF}(\hat{K}', \mathcal{H}(c))$
- 6: **else**
- 7: **return** $K = \text{KDF}(z, \mathcal{H}(c))$
- 8: **end if**

Fig. 2: Description of KYBER.CCAKEM algorithm from [12].

TABLE I: Parameters of CRYSTALS-Kyber.

	n	k	q	η_1	η_2	(d_u, d_v)
KYBER512	256	2	3329	3	2	(10, 4)
KYBER768	256	3	3329	2	2	(10, 4)
KYBER1024	256	4	3329	2	2	(11, 5)

individually and concatenates the output byte arrays. The $\text{Compress}_q(x, d)$ and $\text{Decompress}_q(x, d)$ functions are defined as follows:

$$\text{Compress}_q(x, d) = \lceil (2^d/q) \cdot x \rceil \bmod^+ 2^d,$$

$$\text{Decompress}_q(x, d) = \lceil (q/2^d) \cdot x \rceil.$$

The functions \mathcal{G} and \mathcal{H} represent the SHA3-512 and SHA3-256 hash functions, respectively. The KDF represents key derivation function. It is realized by SHAKE-256.

There are three different parameter sets for CRYSTALS-Kyber, called KYBER512, KYBER768 and KYBER1024, see Table I. In this paper, we focus on KYBER768 with the security level $k = 3$.

IV. EXPERIMENTAL SETUP

This section describes the equipment we use for trace acquisition and the FPGA implementation of CRYSTALS-Kyber from [7].

A. Equipment

Our equipment is shown in Fig. 3. It consists of the ChipWhisperer-Lite board and the CW305 Artix-7 FPGA target board.

The ChipWhisperer toolkits include hardware and open-source software for security analysis based on a low-cost hardware platform [28]. The ChipWhisperer-Lite board is integrated with high-speed power measurement and programmer for the target device. The board is capable of capturing power traces of the target device synchronously, and controls the communication between the target device and the computer. A shunt resistor is placed between the power supply and the target device for power measurement. The ChipWhisperer-Lite board has a maximum sampling rate of 105 MS/sec and a buffer size of 24,400 samples.

The target board CW305 contains an Artix-7 XC7A100T FPGA.

B. Target FPGA implementation

Our target is the FPGA implementation of CRYSTALS-Kyber presented in [7]. It does not contain any countermeasures against side-channel attacks.

The implementation is described in Verilog. Its architecture consists of a client module, which is responsible for encapsulation, and a server module which is responsible for key

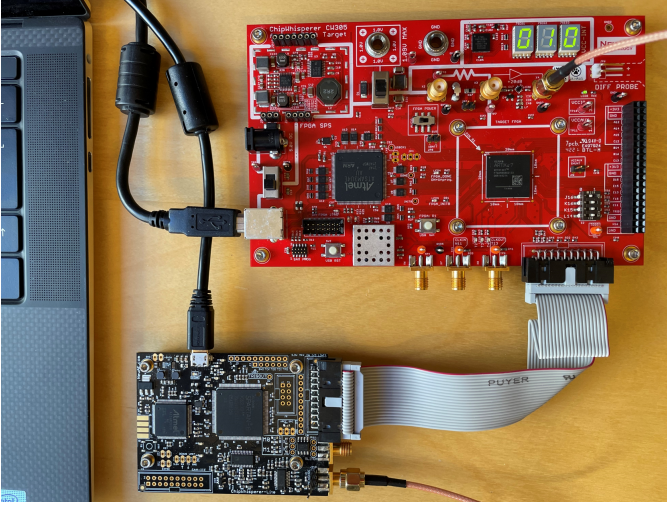


Fig. 3: Equipment for trace acquisition.

generation and decapsulation. The implementation achieves good performance with limited resources. It can fit into the smallest Xilinx Artix-7 device. The implementation covers all three versions of CRYSTALS-Kyber.

We adapted the implementation to a CW305 target board and complemented it with a controller in Verilog which sends in keys and ciphertexts to the server and reads out session keys from the server. We synthesized the Kyber768 implementation to target a Xilinx XC7A100TFTG256 FPGA device using the Vivado synthesis toolchain.

V. PROFILING STAGE

This section describes our profiling strategy.

A. Training trace acquisition

Using the equipment described in the previous section, we captured traces representing the execution of the decapsulation procedure `KYBER.CCAKEM.Decaps()`. For the profiling stage, the traces are captured for ciphertexts generated by `KYBER.CPAPKE.Enc()` for random messages. Since Kyber is a public-key cryptographic algorithm, we can create a properly generated ciphertext for any message using the corresponding public key. Thus, the device under attack can be used for creating a labeled dataset for profiling [17].

Fig. 4(a) shows a trace representing the execution of the complete `KYBER.CCAKEM.Decaps()` procedure.

B. Location of points of interest

Since, for each trace \mathcal{T}_j of the profiling set \mathcal{T} , the value of the message m_j processed by `KYBER.CCAKEM.Decaps()` during the acquisition of \mathcal{T}_j is known, we can use Welch's t-test [29] to analyze the leakage and check if it is possible to locate intervals corresponding to the individual message bytes.

For each byte $i \in \{0, 1, \dots, 31\}$, we partition \mathcal{T} into two sets, \mathcal{T}_0 and \mathcal{T}_1 as:

$$\begin{aligned}\mathcal{T}_0 &= \{\mathcal{T}_j \in \mathcal{T} \mid m_j[i] < 128\} \\ \mathcal{T}_1 &= \{\mathcal{T}_j \in \mathcal{T} \mid m_j[i] > 128\},\end{aligned}$$

where $m_j[i]$ is i th byte of m_j , for all $j \in \{1, 2, \dots, |\mathcal{T}|\}$. The t-test determines if there is a noticeable differences in the means of \mathcal{T}_0 and \mathcal{T}_1 by computing:

$$t = \frac{\mu_0 - \mu_1}{\sqrt{\frac{\sigma_0^2}{n_0} + \frac{\sigma_1^2}{n_1}}},$$

where μ_i , σ_i , and n_i are the mean, standard deviation and cardinality of \mathcal{T}_k for $k \in \{0, 1\}$.

Fig. 4(c) shows the t-test results. One can clearly see 32 groups corresponding to the processing of 32 message bytes. Fig. 4(d) gives a zoomed-in view of 8 bytes. We can see two groups of t-test peaks for every byte, separated by a space. Each byte, with the exception of the ones at the beginning and at the end, overlaps with two adjacent bytes. This overlap could mean that several computations involving different message bytes are carried out in parallel. We will discuss this issue in more details in Section VI.

C. Training set expansion by cut-and-join

From Fig. 4(d) we can see that the shape of t-test peaks for different bytes is similar. Hence, we can increase the size of the profiling set by a factor of 32 (the number of bytes) using the cut-and-join technique from [17] without having to capture 32 times as many traces. For each trace $\mathcal{T}_j \in \mathcal{T}$ and each byte $i \in \{0, 1, \dots, 31\}$, we cut an interval of \mathcal{T}_j including all t-test peaks for the byte i . The union of the resulting intervals gives us the expanded profiling set. In our experiments, we use 64-point intervals.

Fig. 5 shows the t-test results after cut-and-joining profiling traces. One can see that the peaks of different bytes are perfectly aligned. In the figure we show only selected bytes in order not to make it too crowded, but the perfect alignment holds for other bytes as well.

D. Trace pre-processing

As in many other side-channel attacks, we pre-process traces by standardization (or *variance scaling* [30]) to smooth the shifting along y-axis of power measurements.

Given a set of traces \mathcal{T} with elements $\mathcal{T} = (\tau_1, \dots, \tau_{|\mathcal{T}|})$, each trace $\mathcal{T} \in \mathcal{T}$ is standardized to $\mathcal{T}' = (\tau'_1, \dots, \tau'_{|\mathcal{T}|})$ such that, for all $i \in \{1, \dots, |\mathcal{T}|\}$:

$$\tau'_i = \frac{\tau_i - \mu_i}{\sigma_i},$$

where μ_i and σ_i are the mean and standard deviation of the traces of \mathcal{T} at the i th data point.

TABLE II: MLP architecture used for message recovery.

Layer type	Output shape	# Parameters
Batch Normalization 1	64	256
Dense 1	64	4160
Batch Normalization 2	64	256
ReLU 1	64	0
Dense 2	256	16640
Softmax	256	0
Total parameters:	21312	
Trainable parameters:	21056	

E. Network architecture and training parameters

We use multilayer perceptron (MLP) neural networks with the architecture shown in Table II.

The neural networks are trained with a batch size of 1024 for a maximum of 100 epochs using early stopping with patience 20. We use Nadam optimizer with a learning rate of 0.01 and a numerical stability constant $\epsilon = 1e-08$. Categorical cross-entropy is used as a loss function to evaluate the network classification error. 70% of the training set is used for training and 30% is left for validation. Only the model with the highest validation accuracy is saved.

VI. ATTACK STAGE

In this section, we first describe the original multi-bit error injection method proposed in [10] for the attacks on software implementations of lattice-based PKE/KEMs and then explain how we adopted it to recover messages from a hardware implementation. We also show how session keys can be used to derived from the recovered messages.

A. Multi-bit error injection

Let $m = (m[0], m[1], \dots, m[31])$ be a message to be recovered, where $m[i]$ is the i th message byte for $i \in \{0, 1, \dots, 31\}$, and $c = (\mathbf{u}, \mathbf{v})$ be the corresponding ciphertext generated by $\text{KYBER.CPAPKE.Enc}()$ for m . It is known that, by subtracting the center of the integer ring \mathbb{Z}_q from the j th coefficient of \mathbf{v} , the j th bit of m can be flipped without changing other bits [15]. Using this property, the authors of [10] extend the attack set for message recovery from a single trace to 256 traces by injecting all possible multi-bit errors to the message bytes as follows.

For $e \in \{0, 1, 2, \dots, 255\}$, the original ciphertext c is modified to c_e such that $\text{CPAPKE.Dec}()$ decrypts c_e to

$$m_e = (m[0] \oplus e, m[1] \oplus e, \dots, m[31] \oplus e).$$

This means that, all bits of $m[i]$ in which the 8-bit binary expansion of e that has the value “1” are flipped. Then, the prediction for each $m[i] \oplus e$ for all possible 256 values of e is made and the result is XORed with e to get $m[i]$.

In this way, the multi-bit error injection method converts a non-differential side-channel attack to a differential one. This makes good predictions possible, with *not-so-perfect* deep learning models, biased towards some labels, since at least one of the traces in the extended attack set has the ground truth label to which the models are biased towards.

B. Sliced multi-bit error injection

In [10], the multi-bit error injection attack method is applied to a software implementations of Saber in which the message bytes are processed sequentially. However, hardware implementations carry out computations in parallel. From Fig. 4, we can see that the t-test peaks of the adjacent bytes are overlapping. If the same error e is injected into all bytes, as in [10], it may be cancelled out, reducing the attack efficiency.

We propose to overcome this problem by using the following technique which we call *slicing*. Instead of injecting an error e into all message bytes, we inject e into every fourth byte. As a result, the attack set is expanded into a (256×4) -trace set captured for the ciphertexts c_e decrypting to messages m_e whose bytes $m_e[i]$ are defined by:

$$m_e[i] = \begin{cases} m[i] \oplus e & \text{if } i \bmod 4 = x, \\ m[i] & \text{otherwise,} \end{cases} \quad (1)$$

where $x \in \{0, 1, 2, 3\}$.

Clearly, the number of slices may be different from four. For implementations with a higher overlapping between the byte’s computations, a larger number may be more suitable.

Fig. 6 illustrates the difference between error injection in every byte and in every fourth byte. Each rectangular block represents the computation of some byte value. If the same e is injected into every byte, it may be cancelled out due to the overlapping computation of adjacent bytes. However, if e is injected into every fourth byte, the bytes with injected error do not affect each other.

C. Attack trace acquisition

To recover the message m from a properly generated ciphertext c , we collect the attack dataset as follows. For each $x \in \{0, 1, 2, 3\}$, we construct 256 ciphertexts c_e decrypting to messages m_e defined by (1) and capture a set of 256 traces \mathcal{T}_A during the decapsulation of these ciphertexts.

Eight message bytes $m[i]$ such that $(i \bmod 4) = x$ are recovered from \mathcal{T}_A for each x by giving the segments of traces \mathcal{T}_e representing the processing of $m_e[i] = m[i] \oplus e$ as input to the MLP model \mathcal{N} trained at the profiling stage. For each \mathcal{T}_e , the model \mathcal{N} outputs a score vector $S_{i,e} = \mathcal{N}(\mathcal{T}_e)$ in which the value of the l th element, $S_{i,e}[l]$, is the probability that $m_e[i] = l$, for $l, e \in \{0, \dots, 255\}$.

The most likely label for $m[i]$ among 256 candidates is decided as:

$$\tilde{l} = \arg \max_{l \in \{0, 1, \dots, 255\}} \left(\prod_{e=0}^{255} S_{i,e}[l \oplus e] \right).$$

If $\tilde{l} = m[i]$, the classification is successful. The condition $\tilde{l} = m[i]$ can be verified by checking if the rank of $m[i]$ is zero¹.

¹A *rank* of a byte b is the number of other bytes which have the probability higher than b in the score vector.

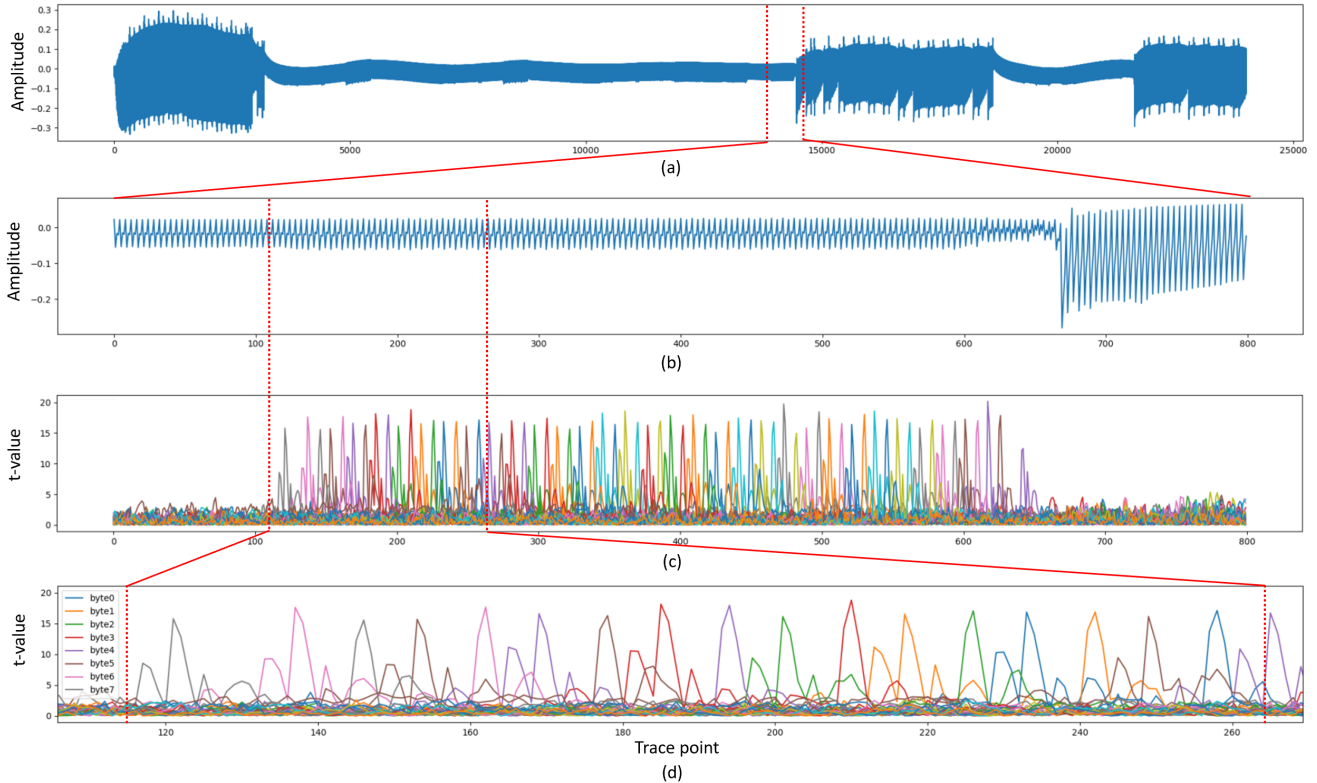


Fig. 4: (a) An average trace representing the execution of `KYBER.CCAKEM.Decaps()`; (b) a segment representing the message processing; (c) t-test results for all 32 message bytes; (d) a zoomed-in view of eight bytes.

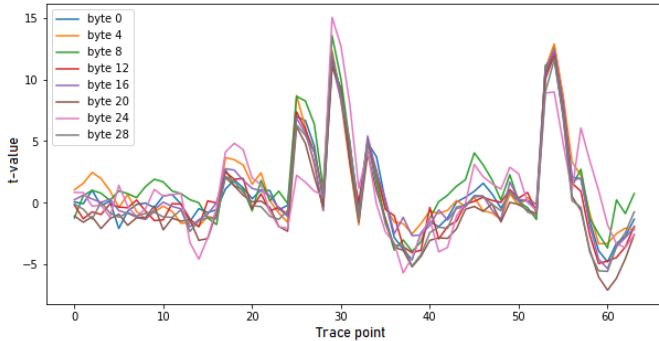


Fig. 5: T-test results for bytes $\{0, 4, \dots, 28\}$ after cut-and-join.

D. Session key recovery

A successful recovery of the message m from a properly generated ciphertext c trivially implies the session key recovery, since the session key can be derived as $K = \text{KDF}(\hat{K}, \mathcal{H}(c))$ where $(\hat{K}, r) = \mathcal{G}(m, \mathcal{H}(pk))$ (see lines 3 and 5 of `KYBER.CCAKEM.Encaps()`).

VII. EXPERIMENTAL RESULTS

This section summarises our experimental results. Following the profiling strategy in Section V, we captured 200K training

traces during the execution of `KYBER.CCAKEM.Decaps()` procedure with input ciphertexts encrypting random messages. The training set was then expanded to 6.4M using the cut-and-join technique. An MLP model with the architecture listed in Table II was trained. In Section VII-B we describe how model analysis helped us to optimize the model.

A. Message/session key recovery attack

In this section, we present the results of message recovery attacks. We compare the proposed sliced multi-bit error injection method to the original multi-bit error injection method and to a direct repetition attack. All attacks are performed on the same 10 messages selected at random. Table III shows the time it takes to capture the test set for one message in each case.

1) *Sliced multi-bit error injection attack*: For each message, we captured 256×4 traces for the 4-sliced multi-bit error injection attack. Each measurement was repeated $N = 5$ times.

Table IV summarizes the results of message recovery. It shows the average and the maximum ranks of all message bytes for each message. We can see that all maximum ranks are smaller than 4. This means that we can recover the messages by enumerating the bytes with ranks 0,1,2 and 3, i.e. using up to $4^{32} = 2^{64}$ enumerations in total, which is feasible.

Fig. 7 shows how the enumeration complexity varies for each message as we increase the number of traces. The y-axis

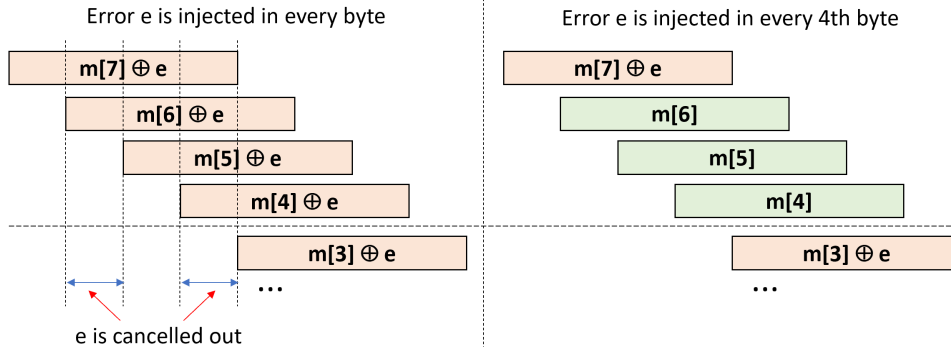


Fig. 6: An illustration of how the injected errors might be cancelled out. The blocks represent computations executed in parallel.

TABLE III: Time for capturing the training and test sets.

	Attack method	# Traces	Time (min)	Success rate
Test set	repetition	256 x 5	5:40	0%
	multi-bit error injection w/o slicing	256 x 5	5:40	10%
	multi-bit error injection w/o slicing	256 x 20	13:21	10%
	4-sliced multi-bit error injection	256 x 4 x 5	22:42	100%
Training set	6 hours to capture 200K traces			

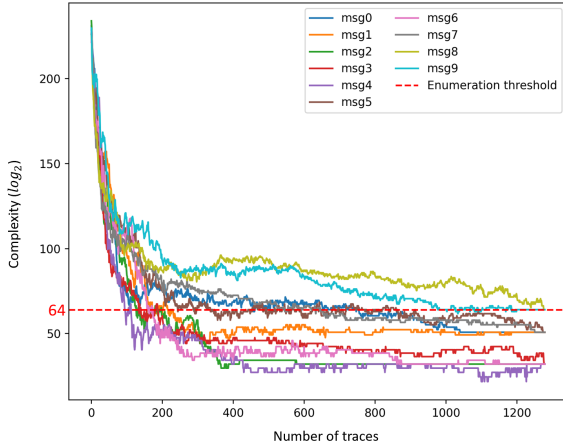


Fig. 7: Enumeration complexity as a function of the number of traces in the 4-sliced multi-bit error injection attack.

represents the number of enumerations required for message recovery in \log_2 scale. The red dashed line marks the threshold $64 = \log_2(2^{64})$ which we consider feasible. We can see that, as the number of traces increases, the enumeration complexity first sharply decreases, and then flattens. Further increase of the number of repetitions N typically does not affect the ranks significantly in our experiments.

2) *Multi-bit error injection attack*: To compare the presented 4-sliced method to the original multi-bit error injection method, we also captured 256 traces in which the error e is

injected in all message bytes. We tested two cases with a different number of repetitions: $N = 5$ and $N = 20$. The case of $N = 20$ is included because the total number of traces used by the original multi-bit error injection method with $N = 20$ is equal to the total the total number of traces used by the 4-sliced multi-bit error injection method with $N = 5$.

Table V summarizes the results of message recovery. We can see that, for both $N = 5$ and $N = 20$ cases, only one message can be recovered with enumeration up to $\leq 2^{64}$, message 2.

Fig. 8 shows how the enumeration complexity changes as the number of traces grows for the case of $N = 5$. We can see that the curves flatten quite quickly. Therefore, further increase of N from from 5 to 20 does not affect the ranks significantly.

This experiment shows that slicing is essential for an attack on the hardware implementation of CRYSTALS-Kyber from [7] in which the message bytes are processed in parallel. Without slicing, only 10% of messages can be recovered with enumeration up to 2^{64} . With slicing, we can recover all messages with the same enumeration.

3) *Repetition attack*: To further highlight the effectiveness of the presented sliced multi-bit error injection method, we also compared it to a direct repetition attack which uses multiple traces captured for the same ciphertext c .

Table VI shows the results of message recovery from 256×5 traces. Clearly, these ranks are not feasibly enumerable. We did not attempt to increase to the number of test traces further because it was obvious that the repetition attack does not work on the target implementation.

TABLE IV: Results of the 4-sliced multi-bit error injection attack using $256 \times 4 \times 5$ traces.

Error injected into every 4th byte	Message									
	0	1	2	3	4	5	6	7	8	9
Avg. byte rank	0.09	0.19	0.09	0.09	0.03	0.19	0.06	0.22	0.28	0.16
Max. byte rank	2	2	1	1	1	2	1	2	3	3

All messages can be recovered with enumeration $\leq 2^{64}$

TABLE V: Results of the multi-bit error injection attack without slicing.

Error injected into all bytes		Message									
		0	1	2	3	4	5	6	7	8	9
256×5 traces	Avg. byte rank	1.63	2.34	0.34	0.84	0.72	1.59	1.93	1.88	2.28	1.09
	Max. byte rank	15	24	2	13	5	23	25	24	26	8
256×20 traces	Avg. byte rank	1.56	1.94	0.31	0.69	0.78	1.22	1.88	1.88	2.16	1.06
	Max. byte rank	16	19	2	11	6	15	20	23	22	9

1 out of 10 messages can be recovered with enumeration $\leq 2^{64}$.

TABLE VI: Results of the repetition attack using 256×5 traces.

Repetition attack	Message									
	0	1	2	3	4	5	6	7	8	9
Avg. byte rank	116	123	127	120	121	99	130	105	138	116
Max. byte rank	251	237	247	248	253	240	255	248	252	255

No messages can be recovered with enumeration $\leq 2^{64}$.

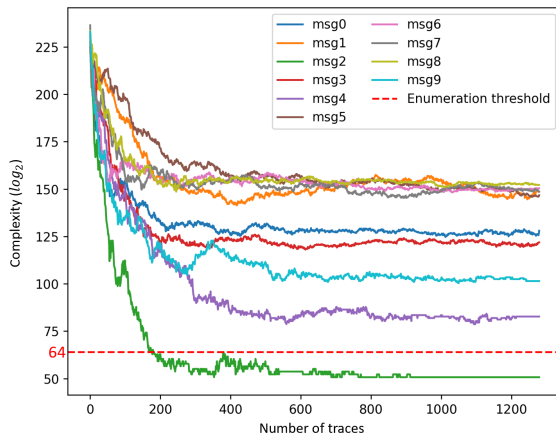


Fig. 8: Enumeration complexity as a function of the number of traces in the multi-bit error injection attack without slicing.

B. Feature and model analysis

Explaining how deep learning models make decisions is important because it may help identify and patch vulnerabilities in the implementation under attack. It may also help in optimizing the model.

1) *Feature analysis*: To evaluate the importance of different input features for the model, we use two techniques: (1) weight analysis, and (2) stuck-at-0 fault injection. Both techniques have been shown useful in previous deep learning-based side-channel attacks of lattice-based PKE/KEMs [31].

Fig. 9 shows a plot representing the weights of the input Batch Normalization layer of the MLP used in the experiments. The Batch Normalization first standardizes the input values X of the layer using their respective mean, μ , and standard deviation, σ , $X_{norm} = (X - \mu)/\sigma$, and then applies the scaling, γ (gamma) and offset, β (beta), parameters to the result, $X' = (\gamma * X_{norm}) + \beta$. The parameters γ and β are learned by the model during the training process, i.e. the backpropagation algorithm is adjusted to operate on the transformed inputs, and error is used to update the new scaling and offset parameters learned by the model. Thus, a higher value of γ indicates the higher importance of the corresponding input point in the decision taken by the model.

We can see that there are two clusters of peaks - one around the point 50 and another - around the point 25. We can relate these clusters to the peaks of the t-test shown in Fig. 5.

The relation becomes even more apparent after the stuck-at-0 fault injection analysis. Fig. 10 shows how the prediction accuracy of the model is affected by setting all but 9 consecutive points of a test trace to 0 before making inference (implying that the model takes its decision based on these 9 points only). If the rank remains close to 0 after the stuck-at-0 faults are injected into all points except for $\{p-4, \dots, p+4\} \pmod{64}$, the point p is important.

In Fig. 10 we can clearly see two “valleys” around the points 50 and 25. They show that the most input features for the model’s decision are contained in these intervals. This, in turn, implies that the computations performed by the implementation of CRYSTALS-Kyber during the corre-

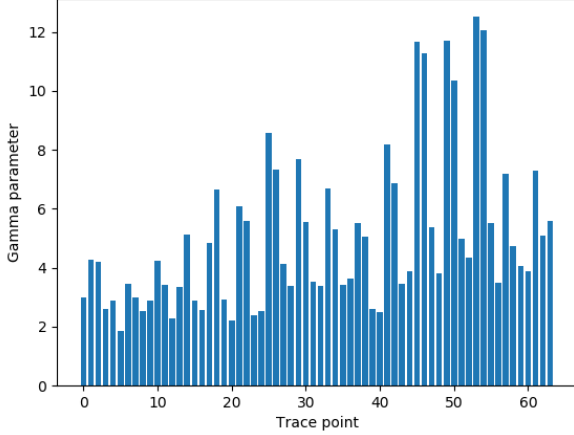


Fig. 9: The γ parameters of the input Batch Normalization layer. Greater γ indicates a more important data point.

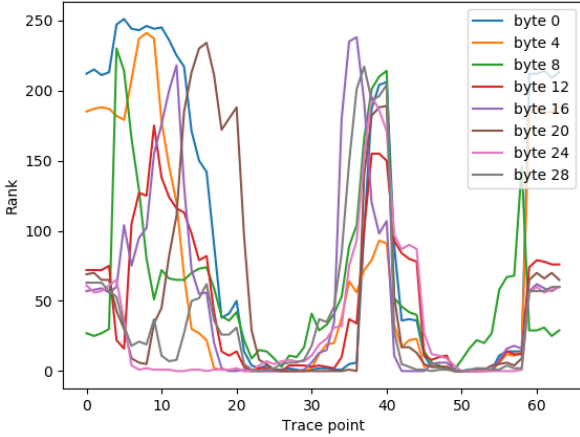


Fig. 10: The effect of stuck-at-0 faults on the rank. For each data point $p \in \{0, \dots, 63\}$, all points except $\{p-4, \dots, p+4\} \pmod{64}$ are set to 0 and inference is performed on the modified data. If the rank remains close to 0, p is important.

sponding clock cycles leak side-channel information. We are currently working on matching the identified points of interest to the CRYSTALS-Kyber procedures in order to re-design the implementation to make it more SCA-resistant.

2) *Model analysis*: In early experiments, we used MLPs with the architecture similar to the one in Table II but with 6 dense layers instead of two. We were expecting models for breaking a hardware implementation of CRYSTALS-Kyber to be more complex than the ones used in the attacks on its software implementations. For example, the MLPs used in [10] for recovering message bytes from an ARM Cortex-M4 CPU implementation of CRYSTALS-Kyber has 3 dense layers.

After multiple training attempts, our best model was able to

TABLE VII: The impact of the number of dense layers on the maximum rank of MLPs. Each dense layer has 256 neurons.

# Layers	Message										Avg	Max
	0	1	2	3	4	5	6	7	8	9		
6	4	2	1	1	1	2	1	2	6	5	2.5	6
5	2	2	1	1	1	2	1	2	9	5	2.6	9
4	2	2	1	1	1	3	1	4	6	2	2.3	6
3	2	2	1	2	1	2	1	3	6	6	2.6	6
2	2	2	1	2	1	2	1	4	6	5	2.6	6

TABLE VIII: The impact of the layer size on the maximum rank of MLPs with two dense layers.

# Neurons	Message										Avg	Max
	0	1	2	3	4	5	6	7	8	9		
256	2	2	1	2	1	2	1	4	6	5	2.6	6
128	2	2	1	2	1	2	1	3	6	6	2.6	6
64	2	2	1	1	1	2	1	2	3	3	1.8	3
32	2	2	1	2	1	2	1	2	6	6	2.5	6
16	2	2	1	2	1	2	1	2	8	7	2.8	8
8	5	2	1	2	2	2	1	3	8	5	3.1	8

recover 70% of messages with the enumeration less than 2^{64} (see the first line of Table VII). We used the 4-slice multi-bit error injection method for the attack, with 5 repetitions, as in the experiments for Table IV.

To check if the model is optimal, we first dropped dense layers one-by-one and re-trained. The results are summarised in Table VII). Surprisingly, the model with two dense layers had approximately the same accuracy as the models with more layers. It could still recover 70% of messages with the enumeration less than 2^{64} and had the maximum rank of 6 (see the last line of Table VII).

We continued verifying optimality of the model with two dense layers by incrementally halving the number of neurons in the layer. For each case, we trained three models using the same training set and hyperparameters. Table VIII shows the results for the best of the three cases. We can see that the model with 64 neurons outperforms the others. We selected this model as optimal for the rest of experiments.

We also tried joining multiple models into an ensemble, but this did not improve the prediction accuracy further. Probably the models make dependent errors. We can see from the columns of Table VIII that some messages are easy and some are difficult for SCA. We do not yet have a confirmed explanation of why this is the case. Traces for all messages are captured under the same conditions, so it is unlikely that measurements are to blame. Since there is an overlapping in message byte processing (see Fig. 4 bottom), different combinations of the bytes processed in parallel sum up to different $8n$ -bit words. It is possible that the difficulty of extracting a given byte from this word by SCA depends on other bytes in the word. A larger set of experiments is required to verify this hypothesis.

VIII. CONCLUSION

We demonstrated a message recovery attack on a hardware implementation of CRYSTALS-Kyber by deep learning-based power analysis. The success of the attack is due to the sliced multi-bit error injection method which we introduced in this paper. It is an extension of the original multi-bit error injection method presented in [10] for side-channel analysis of software implementations of lattice-based PKE/KEMs. Software implementations do not require slicing because they execute instructions sequentially. However, our experimental results show that it is essential for hardware implementations. Without slicing, only 10% of messages can be recovered with enumeration up to 2^{64} . With slicing, we can recover all messages with the same enumeration.

Future work includes re-designing the CRYSTALS-Kyber implementation to make it resistant to side-channel attacks.

IX. ACKNOWLEDGMENTS

This work was supported in part by the Swedish Civil Contingencies Agency (Grant No. 2020-11632), the Swedish Research Council (Grant No. 2018-04482) and the Sweden's Innovation Agency Vinnova (Grant No. 2021-02426)

REFERENCES

- [1] D. Moody, "Status Report on the Third Round of the NIST Post-Quantum Cryptography Standardization Process," *Nistir 8309*, pp. 1–27, 2022, <https://nvlpubs.nist.gov/nistpubs/ir/2022/NIST.IR.8413.pdf>.
- [2] "Announcing the commercial national security algorithm suite 2.0," *National Security Agency, U.S Department of Defense*, Sep 2022, https://media.defense.gov/2022/Sep/07/2003071834/-1/-1/0/CSA_CNSA_2.0_ALGORITHMS_PDF.
- [3] J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, P. Schwabe, G. Seiler, and D. Stehlé, "CRYSTALS-Kyber: a CCA-secure module-lattice-based KEM," in *2018 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2018, pp. 353–367.
- [4] J. W. Bos, M. Gourjon, J. Renes, T. Schneider, and C. van Vredendaal, "Masking Kyber: First-and higher-order implementations," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 173–214, 2021.
- [5] J.-P. D'Anvers, M. V. Beirendonck, and I. Verbauwhede, "Revisiting higher-order masked comparison for lattice-based cryptography: Algorithms and bit-sliced implementations," *Cryptology ePrint Archive, Paper 2022/110*, 2022, <https://eprint.iacr.org/2022/110>.
- [6] D. Heinz, M. J. Kannwischer, G. Land, T. Pöppelmann, P. Schwabe, and D. Sprenkels, "First-order masked Kyber on ARM Cortex-M4," *Cryptology ePrint Archive, Paper 2022/058*, 2022, <https://eprint.iacr.org/2022/058>.
- [7] Y. Xing and S. Li, "A compact hardware implementation of CCA-secure key exchange mechanism CRYSTALS-KYBER on FPGA," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 328–356, 2021.
- [8] P. Ravi, S. Sinha Roy, A. Chattopadhyay, and S. Bhasin, "Generic side-channel attacks on CCA-secure lattice-based PKE and KEMs," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2020, no. 3, p. 307–335, Jun. 2020.
- [9] R. C. Rodriguez, F. Bruguier, E. Valea, and P. Benoit, "Correlation electromagnetic analysis on an fpga implementation of crystals-kyber," *Cryptology ePrint Archive, Paper 2022/1361*, 2022, <https://eprint.iacr.org/2022/1361>. [Online]. Available: <https://eprint.iacr.org/2022/1361>
- [10] R. Wang, K. Ngo, and E. Dubrova, "Making biased DL models work: Message and key recovery attacks on Saber using amplitude-modulated EM emanations," *Cryptology ePrint Archive, Paper 2022/852*, 2022, <https://eprint.iacr.org/2022/852>.
- [11] C. Chen *et al.*, "NTRU algorithm specifications and supporting documentation," 2020, <https://csrc.nist.gov/projects/postquantum-cryptography/round-3-submissions>.
- [12] P. Schwabe *et al.*, "CRYSTALS-Kyber algorithm specifications and supporting documentation," 2020, <https://csrc.nist.gov/projects/postquantum-cryptography/round-3-submissions>.
- [13] J. D'Anvers *et al.*, "Saber algorithm specifications and supporting documentation," 2020, <https://csrc.nist.gov/projects/postquantum-cryptography/round-3-submissions>.
- [14] B.-Y. Sim, J. Kwon, J. Lee, I.-J. Kim, T.-H. Lee, J. Han, H. Yoon, J. Cho, and D.-G. Han, "Single-trace attacks on the message encoding of lattice-based KEMs," *IEEE Access*, vol. 8, pp. 183 175–183 191, 2020.
- [15] P. Ravi, S. Bhasin, S. S. Roy, and A. Chattopadhyay, "On exploiting message leakage in (few) NIST PQC candidates for practical message recovery attacks," *IEEE Transactions on Information Forensics and Security*, vol. 17, pp. 684–699, 2021.
- [16] J. Wang, W. Cao, H. Chen, and H. Li, "Practical side-channel attack on masked message encoding in latticed-based KEM," *Cryptology ePrint Archive, Paper 2022/859*, 2022, <https://eprint.iacr.org/2022/859>.
- [17] K. Ngo, E. Dubrova, Q. Guo, and T. Johansson, "A side-channel attack on a masked IND-CCA secure Saber KEM implementation," *IACR Trans. on Cryptographic Hardware and Embedded Systems*, pp. 676–707, 2021.
- [18] R. Ueno, K. Xagawa, Y. Tanaka, A. Ito, J. Takahashi, and N. Homma, "Curse of re-encryption: A generic power/EM analysis on post-quantum KEMs," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2022, no. 1, p. 296–322, Nov. 2021.
- [19] C. Muijdei, A. Beckers, J. M. B. Mera, A. Karmakar, L. Wouters, and I. Verbauwhede, "Side-channel analysis of lattice-based post-quantum cryptography: Exploiting polynomial multiplication," *Cryptology ePrint Archive, Paper 2022/474*, 2022, <https://eprint.iacr.org/2022/474>.
- [20] Z. Xu, O. M. Pemberton, S. Sinha 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*, pp. 1–1, 2021.
- [21] S. Bhasin, J.-P. D'Anvers, D. Heinz, T. Pöppelmann, and M. V. Beirendonck, "Attacking and defending masked polynomial comparison for lattice-based cryptography," *Cryptology ePrint Archive, Paper 2021/104*, 2021, <https://eprint.iacr.org/2021/104>.
- [22] Y. Huang, M. Huang, Z. Lei, and J. Wu, "A pure hardware implementation of CRYSTALS-KYBER PQC algorithm through resource reuse," *IEICE Electronics Express*, vol. 17, no. 17, pp. 20 200 234–20 200 234, 2020.
- [23] L. Botros, M. J. Kannwischer, and P. Schwabe, "Memory-efficient high-speed implementation of kyber on cortex-m4," in *Progress in Cryptology – AFRICACRYPT 2019*, J. Buchmann, A. Nitaj, and T. Rachidi, Eds. Cham: Springer International Publishing, 2019, pp. 209–228.
- [24] Q. Guo, A. Johansson, and T. Johansson, "A key-recovery side-channel attack on Classic McEliece implementations," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2022, no. 4, p. 800–827, Aug. 2022.
- [25] Y. Ji, "A deep learning based side-channel analysis of an FPGA implementation of Saber," Master's thesis, School of Electrical Engineering and Computer Science, KTH, 2022.
- [26] T. Kamucheka, M. Fahr, T. Teague, A. Nelson, D. Andrews, and M. Huang, "Power-based side channel attack analysis on PQC algorithms," *Cryptology ePrint Archive, Paper 2021/1021*, 2021, <https://eprint.iacr.org/2021/1021>.
- [27] E. Fujisaki and T. Okamoto, "Secure integration of asymmetric and symmetric encryption schemes," in *Annual international cryptology conference*. Springer, 1999, pp. 537–554.
- [28] NewAE Technology Inc., "Chipwhisperer," <https://newae.com/tools/chipwhisperer>.
- [29] B. L. Welch, "The generalization of 'student's' problem when several different population variances are involved," *Biometrika*, vol. 34, no. 1-2, pp. 28–35, 1947.
- [30] A. Zheng and A. Casari, *Feature engineering for machine learning: principles and techniques for data scientists*. O'Reilly Media, Inc., 2018.
- [31] R. Wang, K. Ngo, and E. Dubrova, "Side-channel analysis of Saber KEM using amplitude-modulated EM emanations," in *Proc. of Euromicro DSD/SEAA*, 2022.