

KyberSlash: Exploiting secret-dependent division timings in Kyber implementations

Daniel J. Bernstein^{*†}, Karthikeyan Bhargavan^{‡§}, Shivam Bhasin^{¶**}, Anupam Chattopadhyay^{||**},
Tee Kiah Chia^{**}, Matthias J. Kannwischer^{††}, Franziskus Kiefer[§], Thales Paiva^{‡‡^x^{xi}},
Prasanna Ravi^{||**}, and Goutam Tamvada[§]

^{*}University of Illinois at Chicago, Chicago, IL 60607-7045, USA

[†]Academia Sinica, Taiwan

[‡]Inria, Paris, France

[§]Cryspen, Berlin, Germany

[¶]National Integrated Centre for Evaluation, Nanyang Technological University, Singapore

^{||}College of Computing and Data Science, Nanyang Technological University, Singapore

^{**}Temasek Labs, Nanyang Technological University, Singapore

^{††}Quantum Safe Migration Center, Chelpis Quantum Tech, Taipei, Taiwan

^{‡‡}University of Sao Paulo, Brazil

^xFundep, Brazil

^{xi}CASNAV, Brazil

authorcontact-kyberslash@box.cr.yt.to

28 June 2024

Abstract—This paper presents KyberSlash1 and KyberSlash2 – two timing vulnerabilities in several implementations (including the official reference code) of the Kyber Post-Quantum Key Encapsulation Mechanism, currently undergoing standardization as ML-KEM. We demonstrate the exploitability of both KyberSlash1 and KyberSlash2 on two popular platforms: the Raspberry Pi 2 (Arm Cortex-A7) and the Arm Cortex-M4 microprocessor. Kyber secret keys are reliably recovered within minutes for KyberSlash2 and a few hours for KyberSlash1. We responsibly disclosed these vulnerabilities to maintainers of various libraries and they have swiftly been patched. We present two approaches for detecting and avoiding similar vulnerabilities. First, we patch the dynamic analysis tool Valgrind to allow detection of variable-time instructions operating on secret data, and apply it to more than 1000 implementations of cryptographic primitives in SUPERCOP. We report multiple findings. Second, we propose a more rigid approach to guarantee the absence of variable-time instructions in cryptographic software using formal methods.

1. Introduction

In 2016, the National Institute of Standards and Technology (NIST) launched a global standardization process for Public Key Encryption (PKE), Key Encapsulation Mechanisms (KEM), and Digital Signatures (DS) that can withstand quantum computer attacks, which is widely recognized under the umbrella term “Post-Quantum Cryptography.” After years of evaluation, NIST announced the first set of four algorithms to be standardized in July 2022. Among these, one algorithm is selected for Public Key Encryp-

tion/Key Encapsulation Mechanisms (PKE/KEM) and three algorithms were selected for Digital Signatures. Kyber [1], a KEM based on the Module Learning With Error (MLWE) problem, is being standardized by NIST as ML-KEM in FIPS203 [2]. We expect to soon witness wide-scale adoption of ML-KEM across a wide-spectrum of computing devices, ranging from the high-end general purpose PCs to mobile phone processors all the way until computationally constrained embedded devices.

Since the announcement of the NIST standardization process, Kyber has garnered significant attention regarding its vulnerability to Side-Channel Attacks (SCA) [3]. This concern was a key focus during the NIST PQC standardization process, where the susceptibility of Kyber to SCA and appropriate protection mechanisms were studied by several reported works in literature [4], [5], [6]. Given its anticipated widespread adoption, the safety of Kyber implementations will be even more important in the future.

In this work, we report the discovery of multiple timing vulnerabilities in the official reference implementation of Kyber,¹ as well as several well-known open-source Kyber implementations. Notably, all these implementations are carefully designed to be constant-time at the level of source code, by avoiding secret dependent branches and memory accesses. However, we identify that compilers can introduce timing vulnerabilities through utilization of instructions that execute in variable-time. In particular, we discover these vulnerabilities being caused by certain subroutines that involve divisions by the Kyber prime $q = 3329$ (written as `KYBER_Q` in the code).

1. <https://github.com/pq-crystals/kyber>

1.1. Division variations

It is well known that CPU division instructions are slower for some inputs than for others. It has also been known for a long time that these timing variations might be exploitable; see, e.g., [7, “DIV instruction”] measuring timing variations in the context of fixing Lucky Thirteen.

But this general background does not mean that there is a problem with the divisions in the Kyber reference code. For those divisions, the numerator has a limited range, and the denominator is a compile-time constant (whereas in [7] the denominator was variable). The code is not written to use the CPU’s division instruction, but rather to use divisions in the C programming language. One can reasonably guess that any modern compiler will optimize the division by a constant into a multiplication by a suitable constant; multiplication instructions are well known to be faster than division instructions.

A programmer can easily try an experiment to check this: write the division in C; compile it; check the resulting assembly to see that, yes, the compiler is using a multiplication instruction rather than a division instruction. A pleasant consequence of this automatic optimization is that the resulting binary is unaffected by any potential timing variation from division instructions. (There is still a problem on some embedded processors with variable-time multiplication instructions—see, e.g., [8] and [9]—but that problem is outside the scope of this paper.)

Unfortunately, the experiment described in the previous paragraph isn’t sufficiently systematic, and the guess stated above is an oversimplification. Common changes in compiler options can easily end up producing division instructions instead. For example, asking gcc to optimize for code size (`-Os`) generally disables the conversion of divisions into multiplications. This creates a timing vulnerability. (`-Os` is just one example of the issue: for example, on 32-bit MIPS CPUs, gcc 14.1.0 from May 2024 produces division instructions even when it is optimizing for speed.)

We observe that exactly this vulnerability is triggered by the Kyber reference code. As the vulnerability is caused by the appearance of divisions in Kyber’s C code (`/`), we name the vulnerability KyberSlash. We distinguish two forms of KyberSlash, named KyberSlash1 and KyberSlash2; these arise from different aspects of the cryptography inside Kyber, and turn out to open up very different exploitation mechanisms.

1.2. Our Contribution

The contribution of our paper is manifold.

- We describe two variants of the KyberSlash vulnerability present in the November 2023 version of the Kyber reference implementation: KyberSlash1 is present in the decryption of the CPA-secure encryption scheme underlying Kyber and directly leaks information about the secret key. KyberSlash2 is present in the encryption of the CPA-secure encryption scheme and leaks information about the

ciphertext. While this is unproblematic inside encapsulation as the ciphertext is public, it can be used to construct a plaintext-checking (PC) oracle in decapsulation allowing key recovery.

- We present a practical demo showcasing the exploitability of KyberSlash1 on a Raspberry Pi 2 (Arm Cortex-A7). It crafts special ciphertexts and measures the time for decapsulation running on the same processor in a separate process. It successfully recovers a Kyber512 secret key in 10 out of 10 experiments within 2 to 4 hours.
- We demonstrate the exploitability of KyberSlash2 in a separate demo targeting the Arm Cortex-M4 microcontroller. We craft ciphertexts on a host and send the ciphertexts using serial communication to the target microcontroller which performs a decapsulation and reports back to the host when decapsulation is completed. When timing is performed on the target itself, the attack succeeds for Kyber768 within 4 minutes in 10 out of 10 experiments. Most of this time is spent on transmitting 6144 ciphertexts to the target device. We also demonstrate that the attack still works if timing is performed on a separate attacker device transmitting the ciphertexts to the target device. Consequently, KyberSlash2 is exploitable remotely in certain cases.
- We patch the dynamic analysis tool Valgrind [10] to allow detection of variable-time instructions operating on secret data extending Langley’s ctgrind [11] methodology for detecting timing leaks. With the patched Valgrind, and with modified test programs, we are able to detect the vulnerable division operations in the November 2023 version of the Kyber code. We perform a large scale study with the patched Valgrind and apply it to more than 1000 implementations of various cryptographic primitives within SUPERCOP [12] and identify various potential vulnerabilities due to secret-dependent instruction timings.
- Additionally, we propose a more rigid approach to guarantee the absence of variable-time instructions in cryptographic software by using formal verification.

The code for the two demos is available at <https://kyberslash.cr.yp.to/demos.html>.

1.3. Related work

There is a long literature on side-channel attacks. Attacks often rely on access to physical sensors close to the targeted device; see, e.g., van Eck’s 1985 paper [13] on electromagnetic leaks from monitors, or, as one of many recent examples, consider the EM probe in [14, Section 5]. Modeling and protecting against these information leaks is difficult, with protections continually being broken (see, e.g., [15]) and with security seemingly relying on the hope that attackers are too far away to carry out attacks. Sometimes

attacks exploit a lack of access control for physical sensors *built into* the targeted device, such as the power monitors exploited in [16].

Many other attacks rely purely on timing information. An early example is a timing attack recovering TENEX passwords; see, e.g., [17]. Within cryptography, broad awareness of the power of timing attacks began with Kocher’s 1996 paper [18]. Specific sources of timing variation listed in [18] include “branching and conditional statements” (exploited in that paper), “RAM cache hits” (exploited in [19], [20], and [21]), and “processor instructions (such as multiplication and division) that run in non-fixed time”, along with a general warning about “compiler optimizations” as a source of “unexpected timing variations”. Many timing attacks exploit the fact that attackers are very often allowed to run computations on the target machine (see, e.g., [22]); there have also been some remote timing attacks relying on timing information naturally percolating through networks (see, e.g., [23]). Another avenue for timing attacks comes from the fact that many CPUs vary clock speeds depending on power consumption by default, creating a channel from power monitors to timing; see, e.g., [24] and [25].

By now there are many examples of side-channel attacks against post-quantum cryptography, including systems submitted to the NIST post-quantum competition. For example, [26] targeted non-constant-time error-correcting codes used in LAC, a lattice-based KEM; [27] targeted a non-constant-time ciphertext-comparison operation within FrodoKEM, another lattice-based KEM; and [28] targeted non-constant-time decoding in HQC, a code-based KEM.

Side-channel attacks against KEMs frequently work as follows. The attacker sends maliciously crafted ciphertexts to the decapsulation procedure, such that the decrypted message and its associated variables are related to a targeted portion of the secret key. The attacker uses side channels to obtain information about whether the message decrypted correctly. This reveals incremental information about the secret key, leading to full key recovery once there are enough ciphertexts. Our attack demos follow this pattern but exploit a different side channel, obtaining the first successful timing attacks against the reference implementation of Kyber.

1.4. Disclosure

We have reported KyberSlash to the Kyber team privately shortly after discovery. After discussion with the Kyber team it has been agreed to publicize the vulnerabilities immediately. The main consideration for this was that due to the upcoming standardization later in 2024, the current deployment of the affected code is small, but the potential impact of delaying the announcement would be much more devastating. The divisions have been replaced by explicit multiplications in the official reference implementation prior to the announcement.²

2. KyberSlash1 vulnerability was patched in the commit dda29cc dated December 1, 2023 and KyberSlash2 vulnerability was patched in the commit 272125f dated December 30, 2023.

Starting in December 2023, further patches addressing KyberSlash have been applied to at least the following cryptographic libraries: zig³, kyber-k2so⁴, CIRCL⁵, AWS-LC⁶, Botan⁷, liboqs⁸, crystals-go⁹, PQCclean¹⁰, pqcrypto-kyber¹¹, pypqc¹², pqm4¹³, and kyberlib¹⁴. We have not investigated whether these libraries were exploitable before the patches. There are some other Kyber implementations, such as [29], that never included the problematic divisions.

We have communicated with numerous maintainers of potentially affected libraries and feedback has been uniformly positive.

2. Notation

For any prime q , we denote the field of integers modulo q as \mathbb{Z}_q . When n is a fixed positive integer, we let R_q denote the polynomial ring $\mathbb{Z}_q[x]/(x^n + 1)$. Then, R_q^k is the module of rank k whose scalars are polynomials in R_q . Polynomials $a \in R_q$ are denoted using lowercase letters. Vectors $\mathbf{a} \in R_q^k$ and matrices $\mathbf{A} \in R_q^{k \times k}$ are denoted in bold using lowercase and uppercase, respectively. When $\mathbf{u}, \mathbf{v} \in R_q^k$, we let $\langle \mathbf{u}, \mathbf{v} \rangle \in R_q$ denote their dot product. The i th entry of vector $\mathbf{a} \in R_q^k$ is denoted as $\mathbf{a}[i]$. Similarly, for a polynomial $a \in R_q$, we use $a[i]$ to denote its coefficient associated with the power x^i .

We denote by \mathcal{B}_η the centered binomial distribution (CBD) with range $[-\eta, \eta]$. For a concise notation, we let $\mathbf{a} \leftarrow \mathcal{B}_\eta(R_q^k)$ mean that each coefficient from each polynomial of vector $\mathbf{a} \in R_q^k$ is drawn according to \mathcal{B}_η . Furthermore, we write $\mathbf{a} \leftarrow \mathcal{B}_\eta^r(R_q^k)$ to denote a derandomized sampling where the randomness comes from a string r . Furthermore, $y \leftarrow \text{Compress}(x, d)$ denotes the lossy compression of x to d bits, where $d < \lceil \log_2 q \rceil$. The compression function is defined as $\text{Compress}(x, d) = \lfloor (2^d/q)x \rfloor \bmod 2^d$, where $\lfloor \cdot \rfloor$ denotes the rounding function that rounds up on ties. The decompression is defined as $x' = \text{Decompress}(y, d) = \lfloor (q/2^d)y \rfloor$.

3. Kyber

Kyber is a KEM designed for CCA security based on the Module-Learning With Errors problem (MLWE) [30], [31]. It offers parameter sets designed for NIST security levels 1, 3, and 5. For each security level, it uses fixed parameters $q = 3329$ and $n = 256$ that define the polynomial ring

3. <https://github.com/ziglang/zig>
4. <https://github.com/symbolicsoft/kyber-k2so>
5. <https://github.com/cloudflare/circl>
6. <https://github.com/aws/aws-lc>
7. <https://github.com/randombit/botan>
8. <https://github.com/open-quantum-safe/liboqs>
9. <https://github.com/kudelskisecurity/crystals-go>
10. <https://github.com/PQCclean/PQCclean>
11. <https://github.com/rustpq/pqcrypto>
12. <https://github.com/JamesTheAwesomeDude/pypqc>
13. <https://github.com/mupq/pqm4>
14. <https://github.com/sebastienrousseau/kyberlib>

TABLE 1. KYBER PARAMETERS FOR EACH SECURITY LEVEL [31].

NIST security	Parameter set	k	η_1	η_2	d_u	d_v	Failure probability
Level 1	Kyber512	2	3	2	10	4	2^{-139}
Level 3	Kyber768	3	2	2	10	4	2^{-165}
Level 5	Kyber1024	4	2	2	11	5	2^{-175}

$R_q = \mathbb{Z}_q[x]/(x^n + 1)$, over which most of the operations are performed.

Given a desired security level, the setup takes public parameters k, η_1, η_2, d_u , and d_v from Table 1. Parameter k defines the rank of the modules used in the scheme. Parameters η_1 and η_2 define the centered binomial distributions \mathcal{B}_{η_1} and \mathcal{B}_{η_2} used to generate coefficients with small norm in \mathbb{Z}_q . Integers d_u and d_v are the number of bits into which coefficients from the two parts of the ciphertext are compressed.

Kyber uses an encoding procedure that allows it to recover the message even after some noise accumulates during the encryption and decryption. It encodes a 256-bit message $\mathbf{m} \in \mathbb{Z}_2^{256}$ into a polynomial in R_q as

$$\text{Encode}(\mathbf{m}) = m_0 + m_1x + \dots + m_{n-1}x^{n-1} \in R_q,$$

where $m_i = \mathbf{m}[i] \lceil q/2 \rceil$. In other words, if bit $\mathbf{m}[i] = 0$ then $m_i = 0$, otherwise $m_i = \lceil q/2 \rceil$. A simple decoding procedure can then be applied to a polynomial m' . Namely, the decoding function outputs a 256-bit message $\mathbf{m}' = \text{Decode}(m')$ from a noisy polynomial m' by simply checking if each coefficient of m' is closer to 0 or to $q/2$, modulo q , and decoding it to 0 or 1, correspondingly.

Similar to most proposed post-quantum KEMs, Kyber is built in two layers. The bottom layer consists of a public-key encryption (PKE) scheme that is designed to be secure against passive adversaries, or, more precisely, against chosen-plaintext attacks (CPA). The top layer, which is a key encapsulation mechanism (KEM) designed to be secure against more powerful chosen-ciphertext attacks (CCA), is then constructed by applying a variation of the Fujisaki-Okamoto [32], [33] security conversion over the PKE scheme. Sections 3.1 and 3.2 provide the details of these two layers of algorithms.

3.1. Kyber's auxiliary PKE algorithms designed for CPA security

PKE schemes are defined by three algorithms: key generation, encryption and decryption. These are described by the corresponding procedures listed in Figure 1. The key generation procedure is essentially a creation of an instance of the MLWE problem that protects the secret key. Similarly, the encryption procedure consists of generating another MLWE instance, which now protects the message \mathbf{m} from being recovered from the ciphertext by someone who does not know the secret key \mathbf{sk} .

To see why decryption works, first notice that the ciphertext compression and decompression adds a small noise

```

1: procedure PKE.KEYGEN
2:    $\mathbf{A} \leftarrow$  random element from  $R_q^{k \times k}$ 
3:    $\mathbf{s} \leftarrow \mathcal{B}_{\eta_1}(R_q^k)$ 
4:    $\mathbf{e} \leftarrow \mathcal{B}_{\eta_1}(R_q^k)$ 
5:    $\mathbf{t} \leftarrow \mathbf{A}\mathbf{s} + \mathbf{e}$ 
6:    $\mathbf{pk} \leftarrow (\mathbf{A}, \mathbf{t})$ 
7:    $\mathbf{sk} \leftarrow \mathbf{s}$ 
8:   return  $(\mathbf{pk}, \mathbf{sk})$ 

9: procedure PKE.ENCRYPT( $\mathbf{pk}, \mathbf{m} \in \mathbb{Z}_2^{256}, r$ )
10:   $\mathbf{A}, \mathbf{t} \leftarrow \mathbf{pk}$ 
11:   $\triangleright$ Pseudorandom function PRF is used for sampling
12:   $\mathbf{r} \leftarrow \mathcal{B}_{\eta_1}^{\text{PRF}(r,0)}(R_q^k)$ 
13:   $\mathbf{e}_1 \leftarrow \mathcal{B}_{\eta_2}^{\text{PRF}(r,1)}(R_q^k)$ 
14:   $e_2 \leftarrow \mathcal{B}_{\eta_2}^{\text{PRF}(r,2)}(R_q)$ 
15:   $\mathbf{u} \leftarrow \mathbf{A}^T \mathbf{r} + \mathbf{e}_1$ 
16:   $v \leftarrow \text{Encode}(\mathbf{m}) + \langle \mathbf{t}, \mathbf{r} \rangle + e_2$ 
17:   $\mathbf{c}_u \leftarrow \text{Compress}(\mathbf{u}, d_u)$   $\triangleright$  KyberSlash2
18:   $c_v \leftarrow \text{Compress}(v, d_v)$   $\triangleright$  KyberSlash2
19:  return  $(\mathbf{c}_u, c_v)$ 

20: procedure PKE.DECRYPT( $\mathbf{sk}, (\mathbf{c}_u, c_v)$ )
21:   $\mathbf{u}' \leftarrow \text{Decompress}(\mathbf{c}_u, d_u)$ 
22:   $v' \leftarrow \text{Decompress}(c_v, d_v)$ 
23:   $m' \leftarrow v - \langle \mathbf{u}', \mathbf{s}' \rangle$ 
24:  return  $\text{Decode}(m')$   $\triangleright$  KyberSlash1

```

Figure 1. Kyber's PKE algorithms, with marks indicating where KyberSlash1 and KyberSlash2 appear.

when going from (\mathbf{u}, v) to (\mathbf{u}', v') . Now, by expanding $m' = v - \langle \mathbf{u}', \mathbf{s}' \rangle$, one obtains

$$m' = \text{Encode}(\mathbf{m}) + \langle \mathbf{e}, \mathbf{r} \rangle - \langle \mathbf{s}, \mathbf{e}_1 + \Delta \mathbf{u} \rangle + e_2 + \Delta v,$$

where $\Delta \mathbf{u} = \mathbf{u}' - \mathbf{u}$ and $\Delta v = v' - v$. That is, m' is the sum of the encoded message and a polynomial that is constructed from products and sums of elements whose coefficients came from centered binomial distributions, and are, therefore, small. Kyber security parameters are responsible for ensuring that the coefficients in Δm are small enough so that decryption errors occur only with negligible probability, as shown in Table 1.

3.2. Kyber's KEM algorithms designed for CCA security

The Kyber KEM is defined in Figure 2. The scheme is constructed using an implicit rejection [33] variant of the Fujisaki-Okamoto [32] transform, which is designed to achieve CCA security by combining hash functions H and G with a PKE designed for CPA security. The KEM key generation is essentially the same as its PKE counterpart, except for the fact that a secret string z and the public key \mathbf{pk} are packed into the secret key \mathbf{sk} , to allow for additional verification procedures. The KEM encapsulation takes only \mathbf{pk} and returns a ciphertext \mathbf{c} and a shared key K , that is

```

1: procedure KEM.KEYGEN
2:   pk, skPKE ← PKE.KEYGEN
3:   z ← random 256-bit string
4:   sk ← (skPKE, pk, z)
5:   return (pk, sk)

6: procedure KEM.ENCAPS(pk)
7:   m ← random 256-bit string
8:    $\bar{K}, r \leftarrow G(\mathbf{m}, H(\mathbf{pk}))$ 
9:   c ← PKE.ENCRYPT(pk, m, r)
10:  K ← KDF( $\bar{K}, H(\mathbf{c})$ )
11:  return (c, K)

12: procedure KEM.DECAPS(sk, c)
13:  skPKE, pk, z ← sk
14:  m' ← PKE.DECRYPT(skPKE, c)
15:   $\bar{K}', r' \leftarrow G(\mathbf{m}', H(\mathbf{pk}))$ 
16:  c' ← PKE.ENCRYPT(pk, m', r')
17:  if c = c' then
18:    return K ← KDF( $\bar{K}', H(\mathbf{c})$ )
19:  return K ← KDF(z,  $H(\mathbf{c})$ )

```

Figure 2. Kyber’s KEM algorithms.

computed using a key derivation function (KDF). The main objective of the encapsulation is to make the randomness used for encrypting message \mathbf{m} depend on \mathbf{m} itself by defining r based on $G(\mathbf{m}, H(\mathbf{pk}))$, and sampling the disposable values used for encryption using a cryptographic pseudorandom function PRF. This allows for a quick check, during decapsulation, to see if a ciphertext \mathbf{c} is actually valid or not.

Suppose \mathbf{c} is a chosen ciphertext that was manipulated by the attacker. First, we decrypt \mathbf{c} obtaining \mathbf{m}' , then we reencrypt \mathbf{m}' . Now, even if the ciphertext \mathbf{c} can be successfully decrypted by the PKE algorithm, if we get a different ciphertext after reencrypting \mathbf{m}' using randomness r' defined by $G(\mathbf{m}', H(\mathbf{pk}))$, then we consider \mathbf{c} to an invalid ciphertext. Now, to avoid giving information to an attacker about the validity of the ciphertext they sent, a procedure called implicit rejection is used for deriving the shared key. If the ciphertext was considered valid, then we compute the shared secret based on K' that was derived from \mathbf{m}' . Otherwise, we build a fake shared secret applying the KDF to z and \mathbf{c} . Since the attacker does not know z , the fake shared secret K is indistinguishable from a valid one, thus not revealing additional information, and, since the output is deterministic, repeating the same challenge will result in exactly the same K .

4. KyberSlash

We first start by briefly explaining the adversary model for our attack: The attacker attempts to recover the long-term secret key used by the target’s decapsulation procedure of Kyber. We assume that the attacker has the ability to communicate with the target decapsulation procedure with chosen ciphertexts. This is a standard adversarial model used

in several chosen ciphertext based side-channel attacks [34], [35], [36]. We assume that the attacker is able to observe the execution timing of the decapsulation procedure.

We identified two timing vulnerabilities, which we call KyberSlash1 and KyberSlash2, in implementations of division operations in Kyber. Sections 4.1 and 4.2 explain KyberSlash1 and KyberSlash2 respectively.

4.1. KyberSlash1: Leakage from message decoding

The first timing vulnerability is present within the message decoding operation within the decryption procedure (Line 24 in PKE.Decrypt procedure in Fig. 1). This operation denoted as Decode(m'), essentially converts every coefficient of $m' \in R_q$ into corresponding bit of the decrypted message $\mathbf{m}' \in \mathbb{Z}_2^{256}$. This decoding operation for each coefficient of m' is computed as follows:

$$\mathbf{m}'[i] = (((m'[i] \ll 1) + \text{KYBER_Q}/2)/\text{KYBER_Q}) \& 1;$$

This operation should be implemented in constant time since the message polynomial $m' \in R_q$ is sensitive, and incremental information about m' for chosen ciphertexts can be used to recover the secret key \mathbf{sk} [14], [35], [37]. Refer to Fig 3 for the C code snippet of the message decoding procedure, taken from the official reference implementation of Kyber. Notice that this operation contains a division by the Kyber prime (i.e. KYBER_Q) in Line 11 in Fig. 3. We have added highlighting to our figures to emphasize divisions with secret inputs. We compiled the code using `gcc 14.1` for the x86-64 architecture using the `-Os` compiler optimization flag, instructing `gcc` to optimize for code size. Part of the resulting assembly is shown in Fig. 4 and an `idiv` operation presenting a timing leak is highlighted in red (Line 8). Previous versions of `gcc` result in similar code containing `idiv` instructions. It is important to note that this behavior is not observed for compiler optimization flags `-O0`, `-O1`, `-O2`, `-O3`. We observe similar behavior for many other platforms as shown in Appendix C.

```

1 void poly_tomsg(uint8_t msg[32], const poly *a)
2 {
3   unsigned int i, j;
4   uint16_t t;
5   for(i=0; i<KYBER_N/8; i++) {
6     msg[i] = 0;
7     for(j=0; j<8; j++) {
8       t = a->coeffs[8*i+j];
9       t += ((int16_t)t >> 15) & KYBER_Q;
10      /* Division by Kyber Prime */
11      t = ((t << 1) + KYBER_Q/2)/KYBER_Q & 1;
12      msg[i] |= t << j;
13    }
14  }
15 }

```

Figure 3. C code snippet of message decoding operation, containing the vulnerable division operation by the Kyber prime KYBER_Q .

```

1 ...
2 and ax, 3329
3 add eax, edx
4 movzx eax, ax
5 lea eax, [rax+1664+rax]
6 cdq
7 /* Variable-Time Division Instruction */
8 idiv r10d
9 and eax, 1
10 sal eax, cl
11 inc ecx
12 ...

```

Figure 4. Assembly code snippet of the message decoding operation for a single coefficient, when compiled with `gcc 14.1` for the x86-64 architecture using the `-Os` compiler optimization flag.

4.2. KyberSlash2: Leakage from ciphertext compression

The second timing vulnerability is present within the ciphertext compression operation within the encryption procedure (Line 17-18 in CPA.Encrypt procedure in Fig. 1). The compression procedure essentially compresses each coefficient of the input $u \in R_q$ as follows:

$$\text{Compress}_q(u[i], d) = \lceil (2^d/q) \cdot x \rceil \bmod 2^d;$$

The ciphertext compression operation should also be implemented in constant time, as it leaks information about the recomputed ciphertext within the decapsulation procedure (Line 16 of KEM.Decaps procedure in Fig. 2). The recomputed ciphertext is considered sensitive, and leaks information about the secret key for chosen ciphertexts. Refer to Fig. 5 for the C code snippet for ciphertext compression operation from the official reference implementation of Kyber. This operation contains division by the Kyber prime `KYBER_Q`, similar to that of the message decoding procedure, that is highlighted in red (Line 9). Refer to Fig. 6 for the assembly code snippet of a single iteration of the message decoding operation when compiled with `gcc` and `-Os` where the `idiv` operation is highlighted in red (Line 7). We observe similar divisions for other platforms as shown in Appendix D.

```

1 void poly_compress(uint8_t r[128], const poly *a)
2 {
3     unsigned int i, j; int16_t u; uint8_t t[8];
4     for(i=0; i<KYBER_N/8; i++) {
5         for(j=0; j<8; j++) {
6             u = a->coeffs[8*i+j];
7             u += (u >> 15) & KYBER_Q;
8             /* Division by Kyber Prime */
9             t[j] = (((uint16_t)u << 4) + KYBER_Q/2)
10                /KYBER_Q & 15;
11         }
12         r[0] = t[0] | (t[1] << 4);
13         r[1] = t[2] | (t[3] << 4);
14         r[2] = t[4] | (t[5] << 4);
15         r[3] = t[6] | (t[7] << 4);
16         r += 4;
17     }
18 }

```

Figure 5. C code snippet of the ciphertext compression operation, containing the vulnerable division operation by the Kyber prime `KYBER_Q`.

```

1 ...
2 movzx eax, ax
3 sal eax, 4
4 add eax, 1664
5 cdq
6 /* Variable-Time Division Instruction */
7 idiv r9d
8 and eax, 15
9 mov BYTE PTR [rsp-8+rasi], al
10 ...

```

Figure 6. Assembly code snippet of a single iteration of ciphertext compression operation, when compiled with `gcc 14.1` for the x86-64 architecture using the `-Os` compiler optimization flag.

5. KyberSlash1 demo

This section describes `demo1-pi2`, software demonstrating exploitability of the KyberSlash1 variations in decapsulation timing in the end-of-November-2023 official Kyber512 reference code running under Raspbian (`gcc 8.3.0`) on a Raspberry Pi 2 with a BCM2836 CPU, a quad-core Cortex-A7 running at 900MHz.

5.1. Running the demo

The demo software consists of two files: a script `demo1-pi2.sh` and the main attack code `demo1-pi2.c`. Running three attack experiments is a simple matter of running

```
sh demo1-pi2.sh; sh demo1-pi2.sh; sh demo1-pi2.sh
```

if both files are in the current directory. The script automatically downloads the target Kyber code and compiles and runs the demo. The script assumes that the packages `git`, `time`, and `build-essential` are installed.

A typical experiment takes a few hours. A successful experiment, i.e., an experiment recovering the full Kyber512 secret key, prints `yes`, `eve succeeded`. The demo was observed succeeding ten times in ten experiments, with run times of 2:13:53, 2:04:45, 3:32:45, 1:59:11, 3:23:30, 2:06:31, 2:34:05, 2:04:46, 3:16:21, 2:23:04.

The demo is not guaranteed to succeed: it gives up if it has not found the key from timings of $7 \cdot 2^{18}$ decapsulations. An earlier version of the demo, differing only in the mechanism used to check the computer’s clock, succeeded only twice in three experiments.

5.2. Soft divisions

On this platform, `gcc -Os` converts each division into a call to a division subroutine `divsi3`. The CPU includes a hardware division instruction, but by default `gcc` compiles for an ABI that doesn’t guarantee division instructions.

Checking the cost of the `divsi3` subroutine for divisions of n by 3329, for each n in the range of interest, shows that there is a jump by 20 cycles when the numerator n reaches 3329, a further jump by 2 cycles when n reaches 4096, and a further jump by 1 cycle when n reaches 8192.

5.3. Ciphertext selection

As noted above, this is a demo exploiting KyberSlash1, specifically the division in the line

$$t = ((t \ll 1) + \text{KYBER_Q}/2)/\text{KYBER_Q} \ \& \ 1;$$

applied to each coefficient in the noisy message polynomial $m' = v' - \langle \mathbf{s}, \mathbf{u}' \rangle$, where \mathbf{s} is the secret key and (\mathbf{u}', v') is the decompressed ciphertext.

Each input coefficient t in m' turns into a division of $2t + 1664$ by $q = 3329$. Consequently, there is a big jump in division cost on this platform when t reaches 833, and there are smaller jumps when t reaches 1216 and 3264. The demo chooses ciphertexts (\mathbf{u}, v) so that these timings reveal the coefficients in \mathbf{s} , as explained in the following paragraphs.

A Kyber512 ciphertext (\mathbf{u}, v) consists of three elements $\mathbf{u}[0], \mathbf{u}[1], v$ of $R_q = \mathbb{Z}_q[x]/(x^{256} + 1)$, that is, 256-coefficient polynomials, with each coefficient being an integer modulo $q = 3329$. There are some constraints on the coefficients because ciphertext compression enforces rounding. The secret key \mathbf{s} consists of two elements $\mathbf{s}[0], \mathbf{s}[1]$ of R_q , each coefficient being between -3 and 3 . The polynomial m' mentioned above is $m' = v' - \mathbf{s}[0]\mathbf{u}'[0] - \mathbf{s}[1]\mathbf{u}'[1]$.

Consider what decryption does when $\mathbf{u}'[0] = 72x^{100}$, $\mathbf{u}'[1] = 0$, and $v = 2081 + 2081x + \dots + 2081x^{254} + 208x^{255}$. Note that the final coefficient of v is 208, not 2081. All of the coefficients listed here are compatible with the constraints on compressed ciphertexts.

The coefficient of x^{255} in m' has the form $m'[255] = 208 - 72\mathbf{s}[0][155]$. If $\mathbf{s}[0][155]$ happens to be 3, then $m'[255] \bmod 3329$ is 3321, producing a slow division. Otherwise $m'[255]$ is between 64 and 424, producing a fast division. Other coefficients in m' are between 1865 and 2297, producing slow divisions with no obvious dependence upon secrets.

Consider an optimistic model of total decapsulation time as a constant plus the time for this division. Then $\mathbf{s}[0][155] = 3$ produces slow decapsulation, while other possibilities for $\mathbf{s}[0][155]$ produce fast decapsulation, so decapsulation timings immediately distinguish $\mathbf{s}[0][155] = 3$ from the other possibilities. Replacing 72 with -72 similarly distinguishes -3 from $-2, -1, 0, 1, 2, 3$; replacing 72 with 107, which is another allowed coefficient, distinguishes 2, 3 from $-3, -2, -1, 0, 1$; etc. Replacing $72x^{100}$ with $72x^{101}$ targets $\mathbf{s}[0][154]$ instead of $\mathbf{s}[0][155]$. Exchanging the roles of $\mathbf{u}'[0]$ and $\mathbf{u}'[1]$ targets $\mathbf{s}[1]$ instead of $\mathbf{s}[0]$.

Converting this into a complete attack demo was a conceptually straightforward matter of filtering out the noise that appears in real timings. The details are in Section B. Optimizing this demo turned out to be unimportant, since KyberSlash2 allows a more powerful attack approach, explained in Section 6.

6. KyberSlash2 demo

In this section, we discuss how an attacker can recover the key using leakage from the ciphertext compression step of the reencryption procedure.

6.1. Chosen ciphertexts to extract key information

Let us begin by defining how to build ciphertexts whose decapsulation timings allow the attacker to learn information on the secret key. Each malicious ciphertext is defined by 5 parameters: a 256-bit message $\hat{\mathbf{m}} \in \mathbb{Z}_2^{256}$, followed by four integers $\hat{u}, \hat{v}, \hat{i}$, and \hat{j} . To build the malicious ciphertext, first we compute (\mathbf{u}, v) as

$$v = \text{Encode}(\hat{\mathbf{m}}) + \hat{v}, \text{ and } \mathbf{u}[i] = \begin{cases} -\hat{u}x^{(256-j)}, & \text{if } i = \hat{i}, \\ 0, & \text{otherwise.} \end{cases}$$

Remember that polynomial operations in Kyber are done in $R_q = \mathbb{Z}_q[x]/(x^{256} + 1)$. Then, the malicious ciphertext is the compression of (\mathbf{u}, v) , as defined by Kyber.

To avoid having to deal with additional noise from compression and decompression, it makes sense to choose attacking parameters \hat{u} and \hat{v} that are not significantly affected by the lossy compression. This way, during decapsulation, the malicious ciphertext produces the following noisy message

$$m' = v - \langle \mathbf{s}, \mathbf{u} \rangle = \text{Encode}(\hat{\mathbf{m}}) + \hat{v} + \hat{u}x^{(256-j)}\mathbf{s}[\hat{i}].$$

More specifically, each coefficient of the noisy message m' is defined as

$$m'[j] = \begin{cases} \hat{\mathbf{m}}[0][q/2] + \hat{v} - \hat{u}\mathbf{s}[\hat{i}][j], & \text{if } j = 0, \\ \hat{\mathbf{m}}[j][q/2] - \hat{u}\mathbf{s}[\hat{i}][j], & \text{otherwise.} \end{cases}$$

We can now see that, since $\mathbf{s}[\hat{i}]$ is a polynomial of small coefficients, if \hat{u} is sufficiently small, then $m'[j]$ will be correctly decoded, for all $1 \leq j < 256$. However, whether $m'[0]$ would be correctly decoded depends on how large the noise defined by \hat{v}, \hat{u} and $\mathbf{s}[\hat{i}][j]$ is. Therefore, if an attacker is careful in their selection of \hat{v}, \hat{u} , they may be able to learn the coefficient $\mathbf{s}[\hat{i}][\hat{j}]$ when they have the information on whether m' is correctly decrypted to $\hat{\mathbf{m}}$ or not.

This is the core observation used by what are called plaintext-checking (PC) oracle attacks [4], [34], [35], [36], [38], [39]. PC-oracle attacks are a generic class of attacks that assume access to an oracle that, given a ciphertext \mathbf{c} and a message $\hat{\mathbf{m}}$, returns whether \mathbf{c} was decrypted to $\hat{\mathbf{m}}$ or not. Next we discuss how to build a PC-oracle using the decapsulation time, and how it can be used to recover the secret key.

6.2. PC-oracle attack using decapsulation time

Take a pair of 256-bit messages $\hat{\mathbf{m}}_0$ and $\hat{\mathbf{m}}_1$ that differ only in their first bits, and assume that $\hat{\mathbf{m}}_0[0] = 0$ and $\hat{\mathbf{m}}_1[0] = 1$. Then, it is possible [34], [38] to find a small set of pairs (\hat{u}, \hat{v}) such that the knowledge on whether the malicious ciphertext built with parameters $(\hat{\mathbf{m}}_0, \hat{u}, \hat{v}, \hat{i}, \hat{j})$ is decrypted to $\hat{\mathbf{m}}_0$ or $\hat{\mathbf{m}}_1$ completely characterizes the secret key coefficient $\mathbf{s}[\hat{i}][\hat{j}]$. In our attack, we used the same parameters (\hat{u}, \hat{v}) as the ones used by Ravi et al. [38], which are shown in Table 2.

Now, to use these ciphertexts to mount an attack relying on KyberSlash2, we proceed as follows. Generate a pair of

TABLE 2. THE EFFECT OF (\hat{u}, \hat{v}) IN THE DECODED MESSAGES FOR EACH POSSIBLE SECRET COEFFICIENT, CONSIDERING KYBER768 [38].

$s[\hat{i}][\hat{j}]$	Attack parameters (\hat{u}, \hat{v})			
	(207, 937)	(2, 729)	(106, 521)	(106, -728)
-2	$\hat{\mathbf{m}}_1$	$\hat{\mathbf{m}}_1$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$
-1	$\hat{\mathbf{m}}_1$	$\hat{\mathbf{m}}_1$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$
0	$\hat{\mathbf{m}}_1$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$
1	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$
2	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_0$	$\hat{\mathbf{m}}_1$

messages $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ differing only in their first bits, and let \mathbf{c}_0 and \mathbf{c}_1 be their corresponding uncompressed encryptions using Kyber’s CPA encryption algorithm. Let t_0 and t_1 be the decapsulation times when $\hat{\mathbf{m}}_0$ and $\hat{\mathbf{m}}_1$ are observed after decryption of the malicious ciphertext generated with parameters $(\hat{\mathbf{m}}_0, \hat{u}, \hat{v}, \hat{i}, \hat{j})$. Notice that the randomness used when computing ciphertexts \mathbf{c}_0 and \mathbf{c}_1 come from the hashes of $\hat{\mathbf{m}}_0$ and $\hat{\mathbf{m}}_1$, respectively, thus the encryption of \mathbf{c}_0 and \mathbf{c}_1 should not share any noticeable similarities.

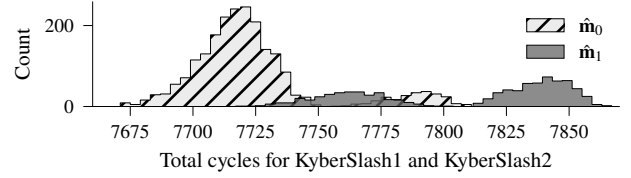
In an idealized scenario where perfect timing is available and KyberSlash2 is the only leakage from the implementation, then t_0 is expected to be slightly different than t_1 . If $t_0 = t_1$, then we can simply restart the process with different pairs $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ until such a difference is observed. This means that, from the decapsulation time, the attacker can infer whether $\hat{\mathbf{m}}_0$ or $\hat{\mathbf{m}}_1$ was observed during decryption. Now, given a pair of indexes (\hat{i}, \hat{j}) , the attacker can verify, for each of the 4 parameters (\hat{u}, \hat{v}) from Table 2, which row matches their observations, thus learning secret coefficient $s[\hat{i}][\hat{j}]$. By iterating over the possible kn index parameters (\hat{i}, \hat{j}) , the attacker then learns the full secret key s .

There are, however, two problems when using this approach in real attacks: in some setups, time measurements come with noise, and the leakages from KyberSlash1 and KyberSlash2 may interfere. In the following, we show how to recover the key in real-world noisy environments.

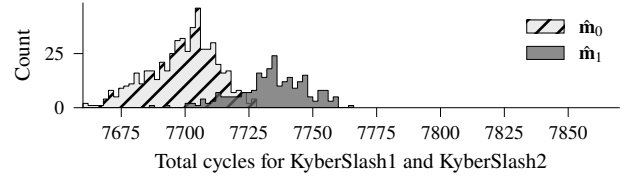
6.3. Key recovery under noisy setups

From the last sections, we know that the problem of key recovery reduces to classifying the observed decapsulation time into $\hat{\mathbf{m}}_0$ or $\hat{\mathbf{m}}_1$. However, because of the interference between KyberSlash1 and KyberSlash2, the distributions that we need to distinguish may not be well separated, as shown in Figure 7a. We now present a series of observations that allow us to simplify this classification by using a careful choice of ciphertexts to be decapsulated.

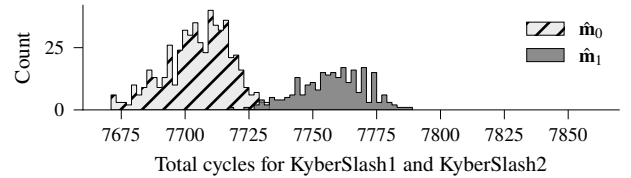
6.3.1. Separating the analysis for pairs (\hat{u}, \hat{v}) . The most important observation is that the leakage from KyberSlash1 is very dependent on the value of \hat{u} . This happens because \hat{u} scales the secret coefficients, resulting in different baselines for the coefficients upon which the message decoding procedure will act. Therefore, we should analyze the leakage distribution for each of the pairs (\hat{u}, \hat{v}) separately. Figure 7b illustrates how, by focusing on only one pair,



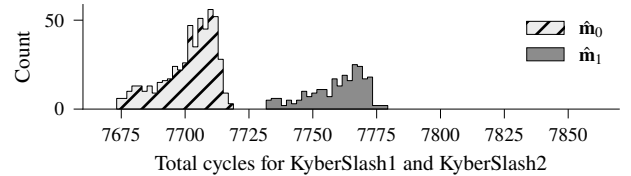
(a) Interference between KyberSlash1 and KyberSlash2.



(b) Leakage within the same group of parameters $(\hat{u}, \hat{v}) = (2, 729)$.



(c) Leakage for $(\hat{u}, \hat{v}) = (2, 729)$ considering the most separated pair $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ of messages among $\pi = 500$ randomly generated pairs.



(d) Average leakage for $(\hat{u}, \hat{v}) = (2, 729)$ using $\mu = 20$ pairs $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$.

Figure 7. The series of observations that allow us to increase the separation of the distributions that we need to distinguish for key recovery. We considered parameters for Kyber768, and the values are simulated using our model for the Cortex-M4 division time.

namely $\hat{u}, \hat{v} = (2, 729)$, we get a simpler pair of distributions to distinguish.

6.3.2. Dealing with random noise. The simplest way to deal with random timing noise is to repeat the measurement a number ρ times and compute some robust measure (e.g. median) over the observed values. However, since the leakage from KyberSlash1 is completely determined by the ciphertext, we cannot get rid of it by repeating the measurement for the same ciphertext. One way of lowering the effect of both measurement and KyberSlash1 noise in the separation is to be more careful in the selection of the pair of messages $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$.

Our idea is to generate a number π of message pairs $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$, and select the pair whose KyberSlash2 leakage is separated by the largest number of cycles. Notice that this can be done offline, using only the target’s public key, and

TABLE 3. CLOCK CYCLES OF `UDIV` INSTRUCTION WITH NUMERATOR n AND DENOMINATOR d ON ARM CORTEX-M4 (STM32F407VG). FOR A SIMPLER DESCRIPTION, WE LET $d_{FL} = 2^{\lfloor \log_2 d \rfloor}$.

Case	Clock cycles	Range of n with $d = 3329$
$d = 0$ or $n = 0$	2	0
$n/d_{FL} < 1$	3	1 to $(2^{11} - 1)$
$n/d_{FL} < 2^4$	5	2^{11} to $(2^{15} - 1)$
$n/d_{FL} < 2^8$	6	2^{15} to $(2^{19} - 1)$
$n/d_{FL} < 2^{12}$	7	2^{19} to $(2^{23} - 1)$
$n/d_{FL} < 2^{16}$	8	2^{23} to $(2^{27} - 1)$
$n/d_{FL} < 2^{20}$	9	2^{27} to $(2^{31} - 1)$
$n/d_{FL} < 2^{24}$	10	2^{31} to $(2^{32} - 1)$
$n/d_{FL} < 2^{28}$	11	–
$n/d_{FL} \geq 2^{28}$	12	–

a model of the division timing for the target device. For the Arm Cortex-M4, the Technical Reference Manual [40] states that a `udiv` instruction takes 2–12 cycles depending on input data. We have reverse engineered the division timings for the common STM32F407VG (present on the STM32F407-Discovery board) and show the results on Table 3. We performed similar (but less extensive) experiments on the STM32L476RG suggesting that the ultra-low-power L4 series has the same division timings. While we present timings for arbitrary denominators, for KyberSlash only the column with $d = 3329$ is relevant showing cross-over points at 1, 2^{11} , 2^{15} , 2^{19} , 2^{23} , 2^{27} , 2^{31} . For $d = 3329$, we have confirmed these timings through exhaustive search over the entire numerator space. This shows that for a fixed denominator, the division timing grows monotonically in the value of the numerator allowing to binary search the cross-over points. For variable denominator, we have picked random denominators and specially formed denominators (very small values and powers of two) and searched for the corresponding cross-over points using binary search. We have performed similar reverse engineering for the more complex application-profile processors Arm Cortex-A55 and Arm Cortex-A72 and present the results in Appendix A.

We can then use the division time model for the target microarchitecture to select the best pair $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ of messages out of π pairs generated at random. Figure 7c shows how the selection among $\pi = 500$ pairs allows to better distinguish between the two messages.

6.3.3. Reducing the noise from KyberSlash1. We observe that it is possible to actively reduce the noise from KyberSlash1 if, instead of using only one message pair $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ for building the malicious ciphertexts, we use a collection of μ pairs. Because each message pair has its own KyberSlash1 leakage baseline, using more pairs and taking the average leakage allow us to significantly reduce the deterministic noise. The effect of averaging the result for $\mu = 20$ pairs is illustrated in Figure 7d.

6.3.4. Key recovery. Synthesizing our 3 methods of dealing with noise, we have the following parameters:

- ρ is the number of repetitions we use for each ciphertext to lower random measurement noise;

- π is the number of candidate pairs $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ we test offline to select only the one whose KyberSlash2 leakage is the most separated. We have used 100,000 for all experiments as it results in sufficient separation and ciphertext generation still terminates within seconds.
- μ is the number of pairs $(\hat{\mathbf{m}}_0, \hat{\mathbf{m}}_1)$ that we actually use when performing the attack.

Therefore, to attack Kyber768 using this setup, the number of decapsulations is $4\rho\mu kn$, where the factor of 4 comes from the number of pairs (\hat{u}, \hat{v}) needed to distinguish the secret coefficients (see Table 2).

Now, given the decapsulation time for each of the $4\rho\mu kn$ ciphertexts, we proceed as follows. First we take the median of the ρ repetitions of each ciphertext and consider this value as the observed time. Then we group each set of ciphertexts with the same parameters $(\hat{u}, \hat{v}, \hat{i}, \hat{j})$, and take the average decapsulation time of the μ resulting values. We are now left with a set of $4kn$ values that we need to classify.

For each pair (\hat{u}, \hat{v}) , we use a Gaussian mixture model (GMM) to give the probability that the decapsulation time is associated with $\hat{\mathbf{m}}_0$ or $\hat{\mathbf{m}}_1$. We remark that the GMM is a rather simple model that does not require any training or complicated parameters. We also tested the performance of k -means unsupervised clustering, but it did not perform as well as the GMM. Furthermore, the probabilities that GMM outputs are very useful when evaluating the likelihood of each row in Table 2 since, not only we are able to find the most suitable value of $\mathbf{s}[i][j]$, we can also rank the different possibilities values by their likelihoods.

6.4. Experiments

We present software demonstrating our attack on Cortex-M4 implementations from pqm4 [41] described in [42]. We target the Kyber768 parameter set, but the attack script is straightforwardly modified to all parameter sets of Kyber. The vulnerable functions (`poly_tomsg`, `poly_compress`, `polyvec_compress`) are verbatim copies of the Kyber reference implementations. We use the 4956a30 version of pqm4 which is the commit before the KyberSlash fixes have been ported to pqm4. We use the Arm GNU compiler toolchain version 13.2.1.Re11 from ¹⁵. We target the STM32F407VG Cortex-M4 (present on the STM32F407-Discovery board).

We implement a simple Python script `m4.py` that can be used to perform the attacks. It takes care of generating the corresponding ciphertexts, assembling the software to be run on the board, flashing the software to the board, executing the experiment, and attempting key recovery. In the end it reports if the secret key was recovered successfully. We present two versions of attack software and describe the experiments and results in more detail in the following.

15. <https://developer.arm.com/downloads/-/arm-gnu-toolchain-downloads>

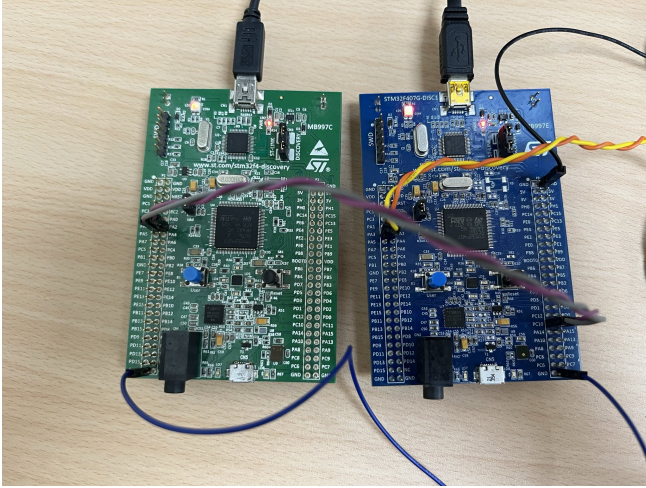


Figure 8. demo2b experiment setup. The target device (left) receives ciphertext from the attacker device (right) through the purple and gray jumper cables (USART3 to USART2) and reports back after decapsulation has completed. The attacker device receives the ciphertexts generated on the host through USART2 (yellow and orange jumper cables), forwards them to the target device through USART3, and reports the timings back to the host.

6.4.1. Local attacker (demo2a). The first version performs cycle counting directly on the target device returning exact cycle counts of the decapsulation to the attacker (the host). Due to the simplicity of the the Cortex-M4 microarchitecture this results in predictable and deterministic cycle counts exactly matching our expectations for timing differences due to KyberSlash1 and KyberSlash2.

As we do not have any noise in this setup, the attack succeeds with a minimum number of measurements: We use the lowest sensible number of messages ($\mu = 2$), and only perform a single measurement ($\rho = 1$) per ciphertext resulting in $n \cdot k \cdot 8 = 6144$ decapsulations. We make use of the 32-bit DWT_CYCCNT cycle counter on the target device allowing us to obtain the exact number of clock cycles a decapsulation took. This cycle count is then sent back to the host using USART which attempts the key recovery. We clock the target at the maximum clock frequency of 168 MHz at which one decapsulation requires approximately 900,000 clock cycles.

An end to end attack takes approximately four minutes out of which only 33 seconds are spent on actual decapsulations, while the remaining overhead is dominated by serial communication with the target device. We achieve significantly higher baud rates using a USB-TTL adapter with a FDTI FT232 chipset rather than the cheaper and more common USB-TTL adapters with a Prolific PL2303 chipset. We determine 806,400 bps as the maximum baudrate for which our setup works reliably.

We perform 10 experiments and successfully recover the secret key in each of them. Running the attack with the above parameters can be done by invoking `m4.py`:

```
./m4.py -i 1 -n 2
```

6.4.2. Remote attacker (demo2b). The second more practical version of the attack software does not make use of

the cycle counter on the target device, but instead performs the timing on the device interacting with the target device. The target reports back using USART after completing the decapsulation, but does not report the cycle counts passed. Note that this adds some noise especially due to the USART clock being much slower than the core clock. For simplicity, we make use of a second Cortex-M4 performing the timing and passing it on to a host laptop. The setup is shown in Figure 8.

The experiments proceed similarly as before, except that there is an intermediate Cortex-M4 that acts as a proxy for the ciphertexts and performs the timing. To improve the reliability of the timing, we make use of interrupts for receiving from the target device (USART3). Between the two boards, we use the highest baud rate that reliably works, which is 768,000 bps for our setup. Furthermore, we clock the target device at a lower clock frequency (24 MHz) than the attacker device (168 MHz) improving the accuracy of our timings.

We achieve reliable key recovery when using $\mu = 2$ and $\rho = 4$. A full experiment takes around 20 minutes out of which 17 minutes are spent on decapsulations. Note that, for $\rho > 1$, we only transmit the ciphertext once which may not be possible for a real attack. One could also increase μ , rather than ρ , however, increasing μ in our setup results in a much more significant increase in runtime due to the slow serial communication. We perform 10 experiments and successfully recover the secret key in each of them. Running the remote attack can be done by executing

```
./m4.py -i 4 -n 2 -r
```

6.4.3. Other platforms. We have performed similar experiments on the Arm Cortex-A55 (which is part of the Qualcomm Snapdragon 888 present in many smartphones) and successfully recovered the secret key when obtaining cycle counts using Perf. However, the measurements in that case contain significantly more noise than in the Cortex-M4 experiments, therefore we require significantly more decapsulations for recovering the key. Nevertheless, in this setup the transmission of the ciphertext is less time consuming than in the M4 experiments, and it is hence preferable to increase μ as discussed before. We have successfully recovered a secret key for $\mu = 40$ and $\rho = 300$ once, but have not yet performed any larger experiments.

7. Detecting secret-dependent divisions

There are many software-analysis tools aimed at systematically ensuring that secrets are not used as branch conditions or as memory addresses, the two most famous sources of timing attacks. The best usability scores in a recent study [43] come from a dynamic-analysis tool, TIME-COP, which has been applied to many existing cryptographic implementations. On the other hand, dynamic-analysis tools inherently catch only what is visible during specific program runs. Those runs do not always achieve the necessary path coverage; the best guarantees require static analysis.

This section investigates what can be done to systematically ensure that secrets are not used as inputs to division instructions. Section 7.1 covers dynamic analysis of many existing cryptographic implementations, and Section 7.2 covers high assurance as a spinoff of formal verification.

7.1. Dynamic scanning using Valgrind

TIMECOP was introduced in [44] as a patch to the SUPERCOP test framework [12]. An evolution of TIMECOP is now maintained as part of SUPERCOP. We have modified TIMECOP to check not just for secret branch conditions and secret load/store addresses but also for secret divisions.

7.1.1. Patching Valgrind. Internally, like various other constant-timeness tools going back to ctgrind [11], TIMECOP runs cryptographic software under the Valgrind tool – specifically, under Valgrind’s Memcheck memory error detector [45]. Our main patch is to Memcheck, and is designed to also be usable via Valgrind-based tools other than TIMECOP, or by developers using Valgrind directly for ad-hoc tests; see Section 7.1.2.

Memcheck applies binary instrumentation to a program to detect, among other things, whether a program uses uninitialized values in non-trivial ways. In particular, it tests for uninitialized branch conditions and uninitialized load/store addresses. Tools such as TIMECOP mark secret data as uninitialized via a Valgrind “client request” [45, §3.1]. Memcheck does not test for uninitialized divisions; this is why a patch is needed.

Our patch, extending a simple 2015 prototype from Dove and Vasiliev [46], issues a warning when a client program uses a division instruction – such as `sdiv` or `udiv` on AArch64, or `idiv` or `div` on x86-64 – to operate on a secret (or uninitialized) data item. The patch also includes preliminary support for catching other variable-latency instructions. We developed the patch for Valgrind 3.22.0 (released October 2023), and it continues to work with Valgrind 3.23.0 (released April 2024).

The patch also modifies Memcheck to print a distinct error message for these operations, to make the operations easier to spot by human readers and scripts. To allow smooth future integration into upstream Valgrind, the patch checks for variable-latency instructions as a run-time option, skipped by default but enabled by a new Valgrind client request `VALGRIND_ENABLE_TIMECOP_MODE`.

7.1.2. Small-scale example: Kyber. We modified the test program `ref/test/test_kyber.c` from the November 2023 Kyber reference code, to invoke the Valgrind client request for the new checking mode, to mark the random number generator’s output as “uninitialized” (i.e. potentially secret), and to mark public key data as “initialized”.

We cross-compiled, for Linux/AArch64, the patched version of Valgrind, as well as the November 2023 version of the Kyber test programs linked with the Kyber512, Kyber768, and Kyber1024 implementations. The Kyber code was built at optimization levels `-O0`, `-O1`, `-Os`, `-O2`, and

```

==7174== Conditional jump or move depends
on uninitialised value(s)
==7174== at 0x108BBC: rej_uniform (indcpa.c:140) ...
==7174== Variable-latency instruction operand
of size 4 is secret/uninitialised
==7174== at 0x1090CC: pqcrystals_kyber512_ref_
polyvec_compress (polyvec.c:48) ...
==7174== Variable-latency instruction operand ...
==7174== at 0x109358: ...poly_compress (poly.c:28) ...
==7174== Variable-latency instruction operand ...
==7174== at 0x10952C: ...poly_tomsg (poly.c:191) ...

```

Figure 9. Sample of Valgrind log showing detection of variable-latency instructions, in modified `test_kyber.c` with Kyber512, compiled with `gcc 11.2.1` for AArch64 with `-Os`

`-O3`, and with debugging information enabled (`-g`). The builds were done on an Apple MacBook Pro (2018) with an Intel x86-64 Core i7, running Alpine Linux 3.19, and with a `gcc 11.2.1` cross-compilation toolchain¹⁶. We then ran the Kyber test programs under Valgrind, using the QEMU¹⁷ emulator.

The Valgrind runs produced instrumentation logs, as partially shown in Figure 9. For the `-Os` binaries, the runs flagged

- lines 28 and 191 of `poly.c` (`poly_compress`, `poly_tomsg`), and line 48 of `polyvec.c`, for Kyber512 and Kyber768,
- and lines 43 and 191 of `poly.c`, and line 24 of `polyvec.c`, for Kyber1024,

as being involved in variable time operations. The flagged instructions correspond to loads of operands for the vulnerable divisions, or operands to be combined with results from the divisions (Figure 10). Moreover, all of the KyberSlash divisions were successfully detected by this method.

We thus show that patching Valgrind can be a practical way to uncover this class of timing vulnerabilities.

7.1.3. Large-scale example: SUPERCOP. Beyond the Valgrind patch, we patched SUPERCOP to provide TIMECOP as part of SUPERCOP’s multi-core dependency-tracking `data-do` tool for collecting and updating a large database of test results, whereas previously SUPERCOP provided TIMECOP only as part of a single-core non-dependency-tracking `do-part` tool aimed at developers testing their own code.

The current development version of SUPERCOP contains 4433 implementations of 1383 cryptographic primitives, all following SUPERCOP’s API, which has also been adopted by various cryptographic competitions and cryptographic libraries. Within these 4433 implementations, 1283 are marked as `goal-constbranch` and `goal-constindex`, meaning that they are designed to avoid secret-dependent branches and array indices. This is also what triggers implementations to be considered by TIMECOP; this does not always mean that they pass TIMECOP.

16. <https://musl.cc>

17. <https://www.qemu.org>

```

...
t[k] = a->vec[i].coeffs[4*j+k];
10cc: 78e27828 ldrsh w8, [x1, x2, lsl #1]
t[k] += ((int16_t)t[k] >> 15) & KYBER_Q;
10d0: 0a887ce2 and w2, w7, w8, asr #31
10d4: 0b080042 add w2, w2, w8
t[k] = (((uint32_t)t[k] << 10)
+ KYBER_Q/2)/ KYBER_Q) & 0x3ff;
10d8: 53163c42 ubfiz w2, w2, #10, #16
10dc: 111a0042 add w2, w2, #0x680
10e0: 1ac70842 udiv w2, w2, w7
...
u = a->coeffs[8*i+j];
1358: 78e27826 ldrsh w6, [x1, x2, lsl #1]
u += (u >> 15) & KYBER_Q;
135c: 0a867ca2 and w2, w5, w6, asr #31
1360: 0b060042 add w2, w2, w6
t[j] = (((uint16_t)u << 4) + KYBER_Q/2)
/ KYBER_Q) & 15;
1364: 531c3c42 ubfiz w2, w2, #4, #16
1368: 111a0042 add w2, w2, #0x680
136c: 1ac50842 udiv w2, w2, w5
...
t = ((t << 1) + KYBER_Q/2)/KYBER_Q) & 1;
msg[i] |= t << j;
152c: 38636805 ldrb w5, [x0, x3]
1530: 531f3c42 ubfiz w2, w2, #1, #16
1534: 111a0042 add w2, w2, #0x680
1538: 1ac60842 udiv w2, w2, w6
...
for(j=0;j<8;j++){
1544: 11000484 add w4, w4, #0x1
msg[i] |= t << j;
1548: 2a050042 orr w2, w2, w5
154c: 38236802 strb w2, [x0, x3]
...

```

Figure 10. Disassembly showing secret operands flagged by patched Valgrind, and corresponding variable-latency instructions, in modified `test_kyber.c` with Kyber512, compiled with `gcc 11.2.1` for AArch64 with `-Os`

For example, two of the primitives are Kyber512 and Kyber768. The Kyber768 primitive has three implementations in SUPERCOP, in three subdirectories `ref`, `compact`, and `avx2` of the `crypto_kem/kyber768` directory. All of these are marked as `goal-constbranch` and `goal-constindex`. The `compact` implementation passes TIMECOP but the `ref` and `avx2` implementations do not. All of the implementations have rejection-sampling loops; the reason the `compact` implementation passes TIMECOP is that it has an extra line of code to declassify the rejection condition.

SUPERCOP compiles each implementation with a list of compilers. The default list includes `gcc -O`, `gcc -O2`, `gcc -O3`, `gcc -Os`, in each case with `-march=native` and `-mtune=native` to optimize for the host CPU, `-fwrapv` to avoid a well-known class of vulnerabilities, and `-fPIC -fPIE` for position independence. The default list also includes five `clang` options. It is important to note that analyzing binaries cannot make any guarantees about what will happen when there are changes in compiler options (e.g., a project not using `-fwrapv`), compiler versions, choice of compiler, and CPU; there is simply the hope that trying more combinations will catch more problems.

We restricted SUPERCOP to the 1283 implementations described above—except that we included SUPERCOP’s

four implementations of New Hope CCA, an ancestor of Kyber, by simply marking them as `goal-constbranch` and `goal-constindex`. As in Section 7.1.2, we added `-g` to the compiler options so that Valgrind output would mention line numbers in source code. We then ran our patched TIMECOP, along with SUPERCOP’s usual tests and benchmarks. We used a dual AMD EPYC 7742 (128 cores in total) with 512GB of RAM. We compiled natively to include, e.g., AVX2 implementations; of course, this also meant that the run was excluding, e.g., Arm implementations. The machine owner had disabled overclocking both for security reasons and for hardware-longevity reasons, so the CPUs were limited to 2.245GHz. The machine is running Debian 12, with `gcc 12.2.0` and `clang 14.0.6`. The run completed in 87 minutes of real time, using 5786 minutes of user time and 193 minutes of system time. Spot-checks during the run showed that all cores were in use at the beginning (with variable RAM usage, typically around 20GB in total for 128 threads), but half the real time was spent waiting for implementations of a few particularly expensive primitives to finish.

This patched TIMECOP run successfully detected various divisions, all of which were specifically the New Hope code with `gcc -Os` (and not `clang -Os`). For example, within `newhope1024cca/avx2`, Valgrind pointed to line 77 of `poly.c`. Manually checking that line finds a division by `NEWHOPE_Q`. Within `newhope1024cca/ref`, Valgrind pointed to lines 16, 41, 82, 83, 84, 85, 115, 116, 354, and 370 of `poly.c`, along with line 215 of `ntt.c`. Manually checking these lines finds that line 16 of `poly.c` is the starting brace of a short function inlined into lines 41, 82, 83, 84, 85, and 115, with a division (actually a mod operator, `%`) on line 19. The other line numbers are directly pointing to divisions in the code.

A separate scan of the New Hope source code finds other division operators, such as an `r / 8` division in `fips202.c`. What distinguishes the TIMECOP results from such a scan is that TIMECOP locates divisions applied to data derived from *secret* inputs.

As a further experiment, we tried adding *all* of the KEMs in SUPERCOP and starting an incremental run. Like the first run, this finished in under 2 hours real time. The output contains 11610 “Variable-latency” lines; the immediately following lines have 2133 different instruction pointers coming from 556 different lines of code in 139 different implementations. A full analysis of those 556 lines of code would be a large project, but here are two illustrative examples. The first report in alphabetical order points to `crypto_kem/hila5/avx2`, specifically a line saying `% (HILA5_Q / 4)` in `kem.c`. The last in alphabetical order points to `crypto_kem/sikep751/ref`, specifically line 263 of `tdiv_qr.c`, which is actually inside the GMP library for big-integer arithmetic. SIKE has been broken in other ways, but this example illustrates the ability of binary analysis to automatically investigate subroutines.

These experiments show that TIMECOP’s data-flow analysis, including our patches, can be efficiently applied to large volumes of existing cryptographic software within

SUPERCOP’s API, in particular producing many examples of data flow from secrets to division instructions. Of course, this does not imply that the examples are exploitable. We also emphasize that the analysis is not a guarantee: it is limited to the binaries created in the TIMECOP run, and to the code paths that are actually taken in the TIMECOP run.

7.2. Using formal methods

Another approach that can help programmers detect and prevent bugs like KyberSlash is the use of formal verification tools. Indeed, KyberSlash1 was first discovered by some of the authors of this paper when they were trying to formally verify a Rust implementation of Kyber.

The security ideal for cryptographic code is *secret independence*, that is, the attacker-observable behavior of a program should not depend on its secret inputs. This means that the program should, of course, not reveal its secrets via its input-output behavior, but also that it does not leak secrets via (say) the program’s runtime or memory accesses.

There are several variations and formal definitions of secret independence defined in the literature that cover different subsets of side-channel attacks. The most common definition seeks to prevent branching on secrets, non-constant-time arithmetic operations (such as division) on secrets, and using secrets as array indexes or memory addresses. This discipline is sometimes called cryptographic constant time. [47]

There are a variety of tools that seek to ensure secret independence in cryptographic code. We refer the reader to recent surveys of these tools and evaluations of their usability for a more complete picture. [43], [48], [49] In the rest of this section, we use the definition of secret independence that is used in the HACL* verified cryptographic library [50], whose formal guarantees are defined and proven for C programs generated from the F* programming language [51].

7.2.1. Secret independence by typechecking. To use any of the formal verification tools on cryptographic code, we must begin by labeling every input and output as either *public* or *secret*. In expressive dependently-typed languages like F*, these secrecy labels can be embedded within the type of each variable, alongside other logical properties needed for correctness (e.g. the range of integers that may be contained in the variable). By default, it is safe to assume that all inputs and outputs within cryptographic code is labeled secret, and the programmer only needs to annotate inputs and outputs that they know to be public.

To verify these labels, we then need to annotate all the primitive operations in the language to reflect our assumptions about whether or not they leak information about their inputs via side-channels. If an operation may leak information about one of its inputs, then that input is labeled as public, preventing cryptographic code from calling it with a secret value. For example, we typically label both inputs to the division operation as public, and on some platforms we may also want to label inputs to certain multiplications as

public. The labels given to language primitives and external libraries are *trusted* and hence must be carefully reviewed to ensure that they capture the operational details of the target platforms.

In the F* library used in HACL*, for example, the types `u8` and `i16` are defined to be *secret integers* and arithmetic operations like division and modulus are not defined for them. Furthermore, secret integers cannot be compared or used as array indexes. All these operations are only available for values declared with the public integer types `pub_u8` or `pub_i16`. Public integers can be converted to secret integers, but converting a secret integer to a public integer requires a call to an explicit `declassify` operation, every use of which needs to be carefully audited.

Given such a secrecy labeling for a program and all the libraries it uses, the type checker can statically verify that the program is secret independent, and point out any parts of the code where the discipline is violated. It is worth noting that such a typing discipline does not really need the full expressiveness of F*; it can easily be implemented in any type system that supports abstract types and interfaces, including Rust and Java.

7.2.2. Finding KyberSlash with F*. The first variant of KyberSlash was found during a larger project of formally verifying a Rust implementation of ML-KEM by translating it to F* and then proving its correctness. As a first step towards a correctness proof, we tried to prove that the translated ML-KEM F* code was secret independent. By using the default integer types, all the inputs and outputs in our code were initially labeled as secret. Then, inputs that are public, such as algorithm parameters or public keys are manually labeled as public by changing their types to use public integers. Finally, outputs that need to be revealed to the application, such as ML-KEM ciphertexts, are *declassified* from secret to public.

When we then run the F* typechecker on the F* code generated from the Rust implementation, it immediately finds and flags the secret dependent division on line 5:

```

1 let compress_q (coefficient_bits: u8) (fe: u16) =
2   let compressed:u32 = (cast (fe <: u16) <: u32) <<!
3     (coefficient_bits +! 1uy <: u8) in
4   let compressed:u32 = compressed +! v_FIELD_MODULUS in
5   let compressed:u32 = compressed !/
6     (v_FIELD_MODULUS <<! 1ul) in
7   get_n_least_significant_bits coefficient_bits compressed

```

To fix this type error, we could label the input field element as `pub_u16` to indicate that it is public, and then prove that this input is indeed public at all call sites, which would fail since this function is used to compress the IND-CPA message coefficient. And if the function was going to be used with secret inputs, we need to rewrite the Rust code to not use division.

Initially, we did both. We wrote a separate function for compressing message coefficients that treated the input field element as secret, and we kept this function for compressing IND-CPA ciphertext coefficients, since we were (incorrectly) assuming that the ciphertexts were public and

declassifying them. Later, when KyberSlash2 was discovered, we fixed our model and moved this declassification from the IND-CPA ciphertexts to the IND-CCA ciphertexts. When we do so, the KyberSlash2 variant also gets flagged by the F* typechecker, and we subsequently reimplemented ciphertext compression as well.

This experience shows both the strengths and the weaknesses of our approach. As long as we correctly annotate and review all public inputs and outputs, the typechecker is able to find secret independence bugs. However, one must be careful when reflecting cryptographic assumptions in the secrecy labeling, or we may miss attacks.

7.2.3. Limitations and Future Directions. When writing code in a compiled language such as Rust or C, the methodology described above ensures that there are no obvious timing leaks. However, even when carefully writing and formally checking the source code, the compiler may produce secret-dependent code.

The formal verification guarantees we obtain from F* above are at the level of source code, and nothing we check here guarantees that the compiler or microarchitecture will not introduce new side-channels that were not visible in the source language semantics.

Modern compilers optimize code aggressively for performance when using high optimization levels. During this process operations such as masking or shifts may be converted into conditional jumps. Techniques as described in 7.1 can be used to analyze the compiled code and detect compiler-introduced secret-dependent operations. To get more formal guarantees, one would need to apply the secret independence checks at the level of machine code, using techniques like those used in the Jasmin assembly implementation of ML-KEM [29].

Acknowledgments

The authors would like to thank Richard Petri and Shih-Ming Yin for their help with setting up the Cortex-M4 demo. Part of this work was funded by FINEP (Financiadora de Estudos e Projetos). Part of this work was funded by the Intel Crypto Frontiers Research Center. Part of this work was funded by the Taiwan’s Executive Yuan Data Safety and Talent Cultivation Project (AS-KPQ-109-DSTCP).

References

- [1] R. Avanzi, J. W. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. Schanck, P. Schwabe, G. Seiler, and D. Stehlé, “CRYSTALS-Kyber (version 3.0): Algorithm specifications and supporting documentation (October 1, 2020),” *Submission to the NIST post-quantum project*, 2020.
- [2] National Institute of Standards and Technology, “FIPS203: Module-lattice-based key-encapsulation mechanism standard (initial public draft),” Federal Inf. Process. Stds. (NIST FIPS), National Institute of Standards and Technology, 2023-08-24 2023, <https://doi.org/10.6028/NIST.FIPS.203.ipd>.
- [3] G. Alagic, D. Apon, D. Cooper, Q. Dang, T. Dang, J. Kelsey, J. Lichtinger, C. Miller, D. Moody, R. Peralta *et al.*, “Status report on the third round of the NIST post-quantum cryptography standardization process,” National Institute of Standards and Technology, Tech. Rep., 2022.
- [4] Y. Tanaka, R. Ueno, K. Xagawa, A. Ito, J. Takahashi, and N. Homma, “Multiple-valued plaintext-checking side-channel attacks on post-quantum KEMs,” *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2023, no. 3, pp. 473–503, 2023.
- [5] P. Ravi, A. Chattopadhyay, J. D’Anvers, and A. Baksi, “Side-channel and fault-injection attacks over lattice-based post-quantum schemes (Kyber, Dilithium): Survey and new results,” *ACM Trans. Embed. Comput. Syst.*, vol. 23, no. 2, pp. 35:1–35:54, 2024. [Online]. Available: <https://doi.org/10.1145/3603170>
- [6] R. Wang, M. Brisfors, and E. Dubrova, “A side-channel attack on a higher-order masked CRYSTALS-Kyber implementation,” in *Applied Cryptography and Network Security - 22nd International Conference, ACNS 2024, Abu Dhabi, United Arab Emirates, March 5-8, 2024, Proceedings, Part III*, ser. Lecture Notes in Computer Science, C. Pöpper and L. Batina, Eds., vol. 14585. Springer, 2024, pp. 301–324. [Online]. Available: https://doi.org/10.1007/978-3-031-54776-8_12
- [7] A. Langley, “Lucky Thirteen attack on TLS CBC,” 2013. [Online]. Available: <https://www.imperialviolet.org/2013/02/04/luckythirteen.html>
- [8] W. de Groot, “A performance study of X25519 on Cortex-M3 and M4,” 2015. [Online]. Available: <https://research.tue.nl/en/studentTheses/a-performance-study-of-x25519-on-cortex-m3-and-m4>
- [9] T. Pornin, “The problem,” 2018. [Online]. Available: <https://www.bearssl.org/ctmul.html>
- [10] N. Nethercote and J. Seward, “Valgrind: a framework for heavyweight dynamic binary instrumentation,” in *28th ACM SIGPLAN Conference on Programming Language Design and Implementation*, June 2007, pp. 89–100, <https://dl.acm.org/doi/10.1145/1250734.1250746>; <https://valgrind.org>.
- [11] A. Langley, “Checking that functions are constant time with valgrind,” 2010, <https://www.imperialviolet.org/2010/04/01/ctgrind.html>.
- [12] D. J. Bernstein and T. Lange, “eBACS: ECRYPT benchmarking of cryptographic systems,” 2024. [Online]. Available: <https://bench.cr.yp.to>
- [13] W. van Eck, “Electromagnetic radiation from video display units: An eavesdropping risk?” *Comput. Secur.*, vol. 4, no. 4, pp. 269–286, 1985. [Online]. Available: [https://doi.org/10.1016/0167-4048\(85\)90046-X](https://doi.org/10.1016/0167-4048(85)90046-X)
- [14] P. Ravi, T. Paiva, D. Jap, J. D’Anvers, and S. Bhasin, “Defeating low-cost countermeasures against side-channel attacks in lattice-based encryption: A case study on Crystals-Kyber,” *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2024, no. 2, pp. 795–818, 2024. [Online]. Available: <https://doi.org/10.46586/tches.v2024.i2.795-818>
- [15] M. Staib and A. Moradi, “Deep learning side-channel collision attack,” *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2023, no. 3, pp. 422–444, 2023. [Online]. Available: <https://doi.org/10.46586/tches.v2023.i3.422-444>
- [16] M. Lipp, A. Kogler, D. F. Oswald, M. Schwarz, C. Eason, C. Canella, and D. Gruss, “PLATYPUS: software-based power side-channel attacks on x86,” in *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*. IEEE, 2021, pp. 355–371. [Online]. Available: <https://doi.org/10.1109/SP40001.2021.00063>
- [17] L. Campbell, “Tenex hackery,” 1991. [Online]. Available: <https://groups.google.com/g/alt.folklore.computers/c/v9KnB8BIXGY/m/aZ-qDLtD0gAJ>

- [18] P. C. Kocher, "Timing attacks on implementations of Diffie-Hellman, RSA, DSS, and other systems," in *Advances in Cryptology - CRYPTO '96, 16th Annual International Cryptology Conference, Santa Barbara, California, USA, August 18-22, 1996, Proceedings*, ser. Lecture Notes in Computer Science, N. Koblitz, Ed., vol. 1109. Springer, 1996, pp. 104–113. [Online]. Available: https://doi.org/10.1007/3-540-68697-5_9
- [19] D. J. Bernstein, "Cache-timing attacks on AES," 2005. [Online]. Available: <https://cr.yp.to/papers.html#cachetiming>
- [20] C. Percival, "Cache missing for fun and profit," 2005. [Online]. Available: <https://www.daemonology.net/papers/htt.pdf>
- [21] E. Tromer, D. A. Osvik, and A. Shamir, "Efficient cache attacks on AES, and countermeasures," *J. Cryptol.*, vol. 23, no. 1, pp. 37–71, 2010. [Online]. Available: <https://doi.org/10.1007/s00145-009-9049-y>
- [22] Z. N. Zhao, A. Morrison, C. W. Fletcher, and J. Torrellas, "Everywhere all at once: Co-location attacks on public cloud FaaS," in *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1, ASPLOS 2024, La Jolla, CA, USA, 27 April 2024- 1 May 2024*, R. Gupta, N. B. Abu-Ghazaleh, M. Musuvathi, and D. Tsafir, Eds. ACM, 2024, pp. 133–149. [Online]. Available: <https://doi.org/10.1145/3617232.3624867>
- [23] D. Brumley and D. Boneh, "Remote timing attacks are practical," in *Proceedings of the 12th USENIX Security Symposium, Washington, D.C., USA, August 4-8, 2003*. USENIX Association, 2003. [Online]. Available: <https://www.usenix.org/conference/12th-usenix-security-symposium/remote-timing-attacks-are-practical>
- [24] Y. Wang, R. Paccagnella, E. T. He, H. Shacham, C. W. Fletcher, and D. Kohlbrenner, "Hertzbleed: Turning power side-channel attacks into remote timing attacks on x86," in *31st USENIX Security Symposium, USENIX Security 2022, Boston, MA, USA, August 10-12, 2022*, K. R. B. Butler and K. Thomas, Eds. USENIX Association, 2022, pp. 679–697. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity22/presentation/wang-yingchen>
- [25] C. Liu, A. Chakraborty, N. Chawla, and N. Roggel, "Frequency throttling side-channel attack," in *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security, CCS 2022, Los Angeles, CA, USA, November 7-11, 2022*, H. Yin, A. Stavrou, C. Cremers, and E. Shi, Eds. ACM, 2022, pp. 1977–1991. [Online]. Available: <https://doi.org/10.1145/3548606.3560682>
- [26] J. D'Anvers, M. Tjepelt, F. Vercauteren, and I. Verbauwhede, "Timing attacks on error correcting codes in post-quantum schemes," in *Proceedings of ACM Workshop on Theory of Implementation Security, TIS@CCS 2019, London, UK, November 11, 2019*, B. Bilgin, S. Petkova-Nikova, and V. Rijmen, Eds. ACM, 2019, pp. 2–9. [Online]. Available: <https://doi.org/10.1145/3338467.3358948>
- [27] Q. Guo, T. Johansson, and A. Nilsson, "A key-recovery timing attack on post-quantum primitives using the Fujisaki-Okamoto transformation and its application on FrodoKEM," in *Advances in Cryptology - CRYPTO 2020 - 40th Annual International Cryptology Conference, CRYPTO 2020, Santa Barbara, CA, USA, August 17-21, 2020, Proceedings, Part II*, ser. Lecture Notes in Computer Science, D. Micciancio and T. Ristenpart, Eds., vol. 12171. Springer, 2020, pp. 359–386. [Online]. Available: https://doi.org/10.1007/978-3-030-56880-1_13
- [28] T. B. Paiva and R. Terada, "A timing attack on the HQC encryption scheme," in *Selected Areas in Cryptography - SAC 2019 - 26th International Conference, Waterloo, ON, Canada, August 12-16, 2019, Revised Selected Papers*, ser. Lecture Notes in Computer Science, K. G. Paterson and D. Stebila, Eds., vol. 11959. Springer, 2019, pp. 551–573. [Online]. Available: https://doi.org/10.1007/978-3-030-38471-5_22
- [29] J. B. Almeida, M. Barbosa, G. Barthe, B. Grégoire, V. Laporte, J. L chenet, T. Oliveira, H. Pacheco, M. Quaresma, P. Schwabe, A. S r , and P. Strub, "Formally verifying Kyber episode IV: implementation correctness," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2023, no. 3, pp. 164–193, 2023.
- [30] A. Langlois and D. Stehl , "Worst-case to average-case reductions for module lattices," *Designs, Codes and Cryptography*, vol. 75, no. 3, pp. 565–599, 2015.
- [31] P. Schwabe, R. Avanzi, J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, G. Seiler, D. Stehle, and J. Ding, "CRYSTALS-KYBER," Technical report, National Institute of Standards and Technology, 2019, <https://csrc.nist.gov/Projects/post-quantum-cryptography/post-quantum-cryptography-standardization/round-3-submissions>.
- [32] E. Fujisaki and T. Okamoto, "Secure integration of asymmetric and symmetric encryption schemes," in *Annual International Cryptology Conference*. Springer, 1999, pp. 537–554.
- [33] D. Hofheinz, K. H velmanns, and E. Kiltz, "A modular analysis of the Fujisaki-Okamoto transformation," in *Theory of Cryptography*, Y. Kalai and L. Reyzin, Eds. Cham: Springer International Publishing, 2017, pp. 341–371.
- [34] G. Rajendran, P. Ravi, J.-P. D'Anvers, S. Bhasin, and A. Chattopadhyay, "Pushing the limits of generic side-channel attacks on LWE-based KEMs-parallel PC oracle attacks on KyberKEM and beyond," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, 2023.
- [35] 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*, pp. 296–322, 2022.
- [36] S. Bhasin, J.-P. D'Anvers, D. Heinz, T. P ppelmann, and M. Van Beirendonck, "Attacking and defending masked polynomial comparison for lattice-based cryptography," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 334–359, 2021.
- [37] P. Ravi, S. S. Roy, A. Chattopadhyay, and S. Bhasin, "Generic side-channel attacks on CCA-secure lattice-based PKE and KEMs," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, vol. 2020, no. 3, pp. 307–335, 2020. [Online]. Available: <https://doi.org/10.13154/tches.v2020.i3.307-335>
- [38] P. Ravi, S. Deb, A. Baksi, A. Chattopadhyay, S. Bhasin, and A. Mendelson, "On threat of hardware trojan to post-quantum lattice-based schemes: a key recovery attack on Saber and beyond," in *International Conference on Security, Privacy, and Applied Cryptography Engineering*. Springer, 2021, pp. 81–103.
- [39] Q. Guo, T. Johansson, and A. Nilsson, "A key-recovery timing attack on post-quantum primitives using the Fujisaki-Okamoto transformation and its application on FrodoKEM," in *Annual International Cryptology Conference*. Springer, 2020, pp. 359–386.
- [40] Arm Limited, "Cortex-M4 Technical Reference Manual." [Online]. Available: <https://developer.arm.com/documentation/ddi0439/latest>
- [41] M. J. Kannwischer, R. Petri, J. Rijneveld, P. Schwabe, and K. Stoffelen, "PQM4: Post-quantum crypto library for the ARM Cortex-M4," <https://github.com/mupq/pqm4>.
- [42] J. Huang, J. Zhang, H. Zhao, Z. Liu, R. C. C. Cheung,  . K. Ko , and D. Chen, "Improved Plantard arithmetic for lattice-based cryptography," vol. 2022, pp. 614–636, Aug. 2022. [Online]. Available: <https://tches.iacr.org/index.php/TCHES/article/view/9833>
- [43] M. Fourn , D. D. A. Braga, J. Jancar, M. Sabt, P. Schwabe, G. Barthe, P.-A. Fouque, and Y. Acar, "These results must be false": A usability evaluation of constant-time analysis tools," 2024, USENIX Security Symposium 2024, to appear. [Online]. Available: <https://www.usenix.org/conference/usenixsecurity24/presentation/fourne>
- [44] M. Neikes, "TIMECOP: automated dynamic analysis for timing side-channels," 2019. [Online]. Available: <https://www.post-apocalyptic-crypto.org/timecop/>
- [45] J. Seward and N. Nethercote, "Using Valgrind to detect undefined value errors with bit-precision," in *USENIX'05 Annual Technical Conference*, April 2005, pp. 17–30, <https://www.usenix.org/legacy/events/usenix05/tech/general/seward.html>.

TABLE 4. CLOCK CYCLES OF `UDIV` INSTRUCTION WITH NUMERATOR n AND DENOMINATOR d ON ARM CORTEX-A55 (SNAPDRAGON 888). FOR A SIMPLER DESCRIPTION, WE LET $d_{\text{FL}} = 2^{\lfloor \log_2 d \rfloor}$.

Case	Clock cycles	Range of n with $d = 3329$
$d = 0$ or $n = 0$	3	0
$n/d_{\text{FL}} < 2^2$	3	1 to $(2^{13} - 1)$
$n/d_{\text{FL}} < 2^6$	4	2^{13} to $(2^{17} - 1)$
$n/d_{\text{FL}} < 2^{10}$	5	2^{17} to $(2^{21} - 1)$
$n/d_{\text{FL}} < 2^{14}$	6	2^{21} to $(2^{25} - 1)$
$n/d_{\text{FL}} < 2^{18}$	7	2^{25} to $(2^{29} - 1)$
$n/d_{\text{FL}} < 2^{22}$	8	2^{29} to $(2^{32} - 1)$
$n/d_{\text{FL}} < 2^{26}$	9	–
$n/d_{\text{FL}} < 2^{30}$	10	–
$n/d_{\text{FL}} \geq 2^{30}$	11	–

- [46] J. Dove and V. Vasiliev, “Automated testing against timing attacks,” Term project, Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, 2015, <https://courses.csail.mit.edu/6.857/2015/projects>.
- [47] G. Barthe, G. Betarte, J. D. Campo, C. Luna, and D. Pichardie, “System-level non-interference of constant-time cryptography. Part II: verified static analysis and stealth memory,” *J. Autom. Reason.*, vol. 64, no. 8, pp. 1685–1729, 2020.
- [48] M. Barbosa, G. Barthe, K. Bhargavan, B. Blanchet, C. Cremers, K. Liao, and B. Parno, “SoK: Computer-aided cryptography,” in *42nd IEEE Symposium on Security and Privacy, SP 2021, San Francisco, CA, USA, 24-27 May 2021*. IEEE, 2021, pp. 777–795.
- [49] J. Jancar, M. Fourné, D. D. A. Braga, M. Sabt, P. Schwabe, G. Barthe, P.-A. Fouque, and Y. Acar, ““They’re not that hard to mitigate”: What cryptographic library developers think about timing attacks,” in *43rd IEEE Symposium on Security and Privacy*. San Francisco: IEEE, 2022.
- [50] J. K. Zinzindohoué, K. Bhargavan, J. Protzenko, and B. Beurdouche, “HACL*: A verified modern cryptographic library.” ACM, 2017, pp. 1789–1806.
- [51] J. Protzenko, J. K. Zinzindohoué, A. Rastogi, T. Ramanandoro, P. Wang, S. Z. Béguélin, A. Delignat-Lavaud, C. Hritcu, K. Bhargavan, C. Fournet, and N. Swamy, “Verified low-level programming embedded in F,” *Proc. ACM Program. Lang.*, vol. 1, no. ICFP, pp. 17:1–17:29, 2017. [Online]. Available: <https://doi.org/10.1145/3110261>

Appendix A. Division leakage for other devices

Table 4 and Table 5 show the reverse engineered division timings for the Arm Cortex-A55 and the Arm Cortex-A72. The timings look similar to the Arm Cortex-M4 timings, but have different cross-over points. The Arm Cortex-A72 has additional shortcuts for power-of-two denominators.

Appendix B. Signal processing for the KyberSlash1 demo

The KyberSlash1 demo takes the following steps to filter out noise in timings, although sufficient noise—or noise dependent on the ciphertexts—can make this fail.

The demo collects batches of measurements. After each batch, it uses all measurements so far to formulate a guess for the Kyber secret key. It recomputes the public key from

TABLE 5. CLOCK CYCLES OF `UDIV` INSTRUCTION WITH NUMERATOR n AND DENOMINATOR d ON ARM CORTEX-A72 (BCM2835). FOR A SIMPLER DESCRIPTION, WE LET $d_{\text{FL}} = 2^{\lfloor \log_2 d \rfloor}$.

Case	Clock cycles	Range of n with $d = 3329$
$d = 0$ or $n = 0$	3	0
power-of-two d	3	–
$n/d < 1$	3	1 to 3328
$n/d_{\text{FL}} < 2^8$	5	3329 to $(2^{19} - 1)$
$n/d_{\text{FL}} < 2^{12}$	6	2^{19} to $(2^{23} - 1)$
$n/d_{\text{FL}} < 2^{16}$	7	2^{23} to $(2^{27} - 1)$
$n/d_{\text{FL}} < 2^{20}$	8	2^{27} to $(2^{31} - 1)$
$n/d_{\text{FL}} < 2^{24}$	9	2^{31} to $(2^{32} - 1)$
$n/d_{\text{FL}} < 2^{28}$	10	–
$n/d_{\text{FL}} \geq 2^{28}$	11	–

this guess, checks for a match, and stops in case of success. It gives up if it has not succeeded after 512 batches.

Each batch includes 7 choices of ciphertexts (\mathbf{u}, \mathbf{v}) targeting each of the 512 positions in the secret. These $7 \cdot 512$ ciphertexts are handled in a random order to limit any impacts from hysteresis. Each ciphertext is tried 16 times in succession, with the timings sorted and only one intermediate timing recorded to remove outliers.

After the batch, the demo computes an interquartile mean of the recorded timings (across all batches) for each of the 7 choices of (\mathbf{u}, \mathbf{v}) at each position, obtaining $7 \cdot 512$ interquartile means. At each position, the demo then obtains a guess for this position of the secret key by comparing

- 7 interquartile means t_0, \dots, t_6 from the *observed* timings and
- for each possibility for this position of the secret key: a *model* of the division timings d_0, \dots, d_6 for the same 7 choices of (\mathbf{u}, \mathbf{v}) .

Specifically, the demo takes the possibility that minimizes the variance of the 7 numbers $t_0 - d_0, \dots, t_6 - d_6$. The point here is that, in the optimistic model from Section 5.3, a correct guess would have these 7 numbers being a constant (and thus variance 0), namely the time for the rest of the decapsulation process. (This is not exactly the same as asking which guess has a model best correlated with the observed timings; it is asking specifically for a *diagonal* correlation. Correlation allows a scaling factor, whereas the model predicts that the scaling factor is 1.)

One expects random noise to settle down at *most* positions before it settles down at *all* positions; a “coupon collector” spends considerable time waiting for the last few coupons. To address this, the demo actually tries multiple guesses after each batch: specifically, it picks 10 positions with the smallest ratios between the second-smallest variance and the smallest variance, and then tries all 2^{10} combinations of first or second guesses at each position.

B.1. Optimizations

No claims of optimality are made for the number of decapsulations used in this demo. Only one timing is kept from each series of 16 timings; presumably the other timings could be used productively, and there is no reason to think

that 16 is optimal. The 2^{10} guesses after each batch could be replaced by more guesses and more advanced lattice attacks; it would be interesting to explore how to optimize “soft-decision decoding” in the lattice context, accounting for varying confidence levels at each position. A ciphertext can target many positions at once, say a random selection of about a third of the 256 positions, and attribute the resulting timing to each of those positions; each position receives some noise from the other positions, but a simple model suggests that this will be outweighed by the speedup.

Appendix C. Assembly Code Snippets for Message Decoding Operation

Fig.11 and Fig. 12 show divisions when compiling Fig 3 for 64-bit and 32-bit Arm processors.

```

1 ...
2 ldrsh w6, [x1, x2, lsl 1]
3 ldrh w2, [x1, x2, lsl 1]
4 and w6, w7, w6, asr 31
5 add w2, w2, w6
6 ubfiz w2, w2, 1, 16
7 add w2, w2, 1664
8 /* Variable-Time Division Operation */
9 udiv w2, w2, w7
10 and w2, w2, 1
11 lsl w2, w2, w4
12 ...

```

Figure 11. Assembly code snippet of the message decoding operation for a single coefficient, when compiled with arm64 gcc 14.1.0 for the AArch64 architecture using the -Os compiler optimization flag.

```

1 ...
2 uxtah r3, lr, r3
3 uxth r3, r3
4 lsls r3, r3, #1
5 add r3, r3, #1664
6 /* Variable-Time Division Instruction */
7 udiv r3, r3, r5
8 and r3, r3, #1
9 lsls r3, r3, r4
10 orr r3, ip, r3
11 ...

```

Figure 12. Assembly code snippet of the message decoding operation for a single coefficient, when compiled with arm-none-eabi-gcc 14.1 for Arm Cortex-M4 CPU (-mcpu=cortex-m4) using the -Os compiler optimization flag.

Appendix D. Assembly Code Snippets for Ciphertext Compression Operation

Fig.13 and Fig. 14 show divisions when compiling Fig 5 for 64-bit and 32-bit Arm processors.

```

1 ...
2 ldrsh w6, [x1, x2, lsl 1]
3 and w2, w5, w6, asr 31
4 add w2, w2, w6
5 ubfiz w2, w2, 4, 16
6 add w2, w2, 1664
7 /* Variable-Time Division Operation */
8 udiv w2, w2, w5
9 and w2, w2, 15
10 strb w2, [x3, x7]
11 add x3, x3, 1
12 cmp x3, 8
13 ...

```

Figure 13. Assembly code snippet of a single iteration of ciphertext compression operation, when compiled with arm64 gcc 14.1.0 for the AArch64 architecture using the -Os compiler optimization flag.

```

1 ...
2 uxth r3, r3
3 lsls r3, r3, #4
4 add r3, r3, #1664
5 cmp r5, r1
6 /* Variable-Time Division Instruction */
7 udiv r3, r3, r4
8 and r3, r3, #15
9 strb r3, [r6], #1
10 bne .L3
11 ..

```

Figure 14. Assembly code snippet of a single iteration of ciphertext compression operation, when compiled with arm-none-eabi-gcc 14.1 for Arm Cortex-M4 CPU (-mcpu=cortex-m4) using the -Os compiler optimization flag.