

Bounding the cache-side-channel leakage of lattice-based signature schemes using program semantics

Nina Bindel, Johannes Buchmann, Juliane Krämer,
Heiko Mantel, Johannes Schickel, and Alexandra Weber

Computer Science Department, TU Darmstadt, Germany
{nbindel, buchmann, jkraemer}@cdc.informatik.tu-darmstadt.de,
{mantel, schickel, weber}@mais.informatik.tu-darmstadt.de

Abstract. In contrast to classical signature schemes, such as RSA or ECDSA signatures, the lattice-based signature scheme ring-TESLA is expected to be resistant even against quantum adversaries. Due to a recent key recovery from a lattice-based implementation, it becomes clear that cache side channels are a serious threat for lattice-based implementations. In this article, we analyze an existing implementation of ring-TESLA against cache side channels. To reduce the effort for manual code inspection, we selectively employ automated program analysis. The leakage bounds we compute with program analysis are sound overapproximations of cache-side-channel leakage. We detect four cache-side-channel vulnerabilities in the implementation of ring-TESLA. Since two vulnerabilities occur in implementations of techniques common to lattice-based schemes, they are also interesting beyond ring-TESLA. Finally, we show how the detected vulnerabilities can be mitigated effectively.

1 Introduction

The threat posed by quantum computers to current public-key cryptography is known since Shor presented a quantum algorithm to solve the factorization and the discrete logarithm problem in polynomial time [31]. How serious this threat is taken became clear, e.g., when NIST announced to start a standardization process for quantum-resistant schemes beginning in fall 2017 [23].

A promising quantum-resistant, also called post-quantum, candidate to substitute current public-key cryptography is lattice-based cryptography that enjoys, among other things, strong security guarantees. However, security guarantees can be undermined by side-channel vulnerabilities at the implementation level. So far, this has happened mostly for implementations of classical cryptography, e.g., [2, 8, 26, 34]. However, first approaches to analyze lattice-based cryptography with respect to side-channel attacks are already made, e.g., [25, 27–29]. Recently, Groot Bruinderink et al. presented the first attack against a lattice-based signature scheme and broke the scheme BLISS [15] using cache side channels of the Gaussian sampling during the signature generation [17]. Although none of the existing lattice-based signature schemes (or their implementations) claim to be secure against side channels, the attack in [17] raises the question how lattice-based signature schemes can be implemented without cache-side-channel

leakage. Furthermore, in the light of NIST’s standardization process, it is important to analyze lattice-based implementations against cache side channels. The scheme ring-TESLA [3], which is one of the most efficient lattice-based signature schemes, seems to be a good candidate to be implemented without cache side channels: During signature generation, ring-TESLA does not use Gaussian sampling, the sampling method that was exploited in the attack on BLISS in [17].

In this work, we use program analysis to compute upper bounds on the cache-side-channel leakage of lattice-based implementations at the example of the signature generation in ring-TESLA¹. More concretely, we follow an approach based on information theory and reachability analysis, which is implemented in CacheAudit [13]. Variants of CacheAudit were used to analyze multiple cryptographic implementations. CacheAudit 0.2 [13] was used to analyze PolarSSL AES and the eSTREAM Profile 1 portfolio (HC-128, Rabbit, Salsa20, and Sosemanuk). CacheAudit 0.2b was used in a systematic study of AES implementations [21]. Another extension of CacheAudit 0.2 was used on modular exponentiation from the libraries libcrypto and OpenSSL [12]. In this article, we extend CacheAudit 0.2b to CacheAudit 0.2c and apply it to ring-TESLA. This is the first analysis of a post-quantum scheme using CacheAudit.

With CacheAudit 0.2c, we determine upper bounds on the leakage of a ring-TESLA implementation for four attacker models. The bounds are sound, i.e., conservative with respect to the attacker models. We obtain upper bounds between 2.6bit and 51.6bit of potential cache-side-channel leakage. By inspecting the code manually, we then identify vulnerable subroutines. We implement countermeasures in the vulnerable subroutines to mitigate the cache-side-channel leakage. Finally, we argue for the effectiveness of the mitigations. For two subroutines, the argument is completely automated by an analysis with CacheAudit 0.2c that reports 0bit leakage. According to our code inspection, a potential for leakage in the signature computation remains, which is intrinsic to a method called rejection sampling. Rejection sampling is used by design in the most efficient lattice-based signature schemes, such as [3, 6, 7, 15]. We argue that the attacker cannot exploit this potential leakage to get information about the secret key. Therefore, we consider our resulting implementation of ring-TESLA to be resistant to the four types of cache-side-channel attacks we consider.

In summary, the contributions of this article are the following.

- We detect four cache-side-channel vulnerabilities in an existing implementation of ring-TESLA by code inspection, selectively supported by automatic program analysis.
- To mitigate the detected vulnerabilities, we augment the ring-TESLA implementation by side-channel countermeasures. We argue for the effectiveness of the countermeasures, again supported by selective program analysis.

¹ We analyze an implementation of ring-TESLA despite an error that was detected in its security reduction, since we expect that reductions given for its predecessor TESLA [4, 6] will be applicable to ring-TESLA as well. Hence, we consider it to be a good candidate for practical applications that require post-quantum signatures.

- To automate parts of our analysis of the unmitigated and mitigated ring-TESLA implementation, we extend the analysis tool CacheAudit 0.2b to CacheAudit 0.2c. More concretely, we implement support for ten additional x86 instructions. The support can be used to analyze occurrences of these instructions in x86 binaries and is not limited to ring-TESLA.

The detection and mitigation of vulnerabilities not only hardens the ring-TESLA implementation against side-channel attacks. Multiple lattice-based primitives, such as key exchange protocols, encryption, and signature schemes, use techniques similar to the ones we analyze in ring-TESLA. In particular, rejection sampling and sparse multiplication, where we find two of the potential vulnerabilities, occur in ring-TESLA, as well as in other lattice-based primitives. Hence, our results also pave the way to make other lattice-based implementations more trustworthy using program analysis.

2 Preliminaries

2.1 Notation

For an integer $n \in \mathbb{N}$, we define $q \in \mathbb{N}$ to be a prime with $q = 1 \pmod{2n}$. We denote the finite field $\mathbb{Z}/q\mathbb{Z}$ with representatives in $[-q/2, q/2] \cap \mathbb{Z}$ by \mathbb{Z}_q . Furthermore, we define $\mathcal{R}_q = \mathbb{Z}_q[x]/\langle x^n + 1 \rangle$ and $\mathcal{R}_{q,[B]} = \{\sum_{i=0}^{n-1} a_i x^i \mid a_i \in [-B, B] \cap \mathbb{Z}\}$ for $B \in [0, q/2] \cap \mathbb{Z}$ and $\mathbb{B}_{n,\omega} = \{\sum_{i=0}^{n-1} a_i x^i \mid a_i \in \{-1, 0, 1\}, \sum_{i=0}^{n-1} |a_i| = \omega\}$ for $\omega \in [0, n] \cap \mathbb{Z}$. All logarithms are in base 2. Let $\sigma \in \mathbb{R}_{>0}$. Let v be a polynomial, then $v \leftarrow_{\sigma} \mathcal{R}$ means sampling each coefficient of v with discrete Gaussian distribution with standard deviation σ and mean 0 over \mathbb{Z} . For a finite set S , we write $s \leftarrow_{\S} S$ to indicate that an element s is sampled uniformly at random from S .

2.2 Description of ring-TESLA

The signature scheme ring-TESLA is parametrized by $n, \omega, d, B, q, U, L, \kappa, \sigma$, by the hash function $H : \{0, 1\}^* \rightarrow \{0, 1\}^{\kappa}$, and by the encoding function $F : \{0, 1\}^{\kappa} \rightarrow \mathbb{B}_{n,\omega}$, see Figure 1. For detailed information about the system parameters and the encoding function F , we refer to the original work [3]. For $c \in \mathbb{Z}$, we denote by $[c]_L$ the unique representative of c in $(-2^{d-1}, 2^{d-1}] \cap \mathbb{Z}$ such that $c = [c]_L$ modulo 2^d and define $[\cdot]_M : \mathbb{Z} \rightarrow \mathbb{Z}, c \mapsto (c - [c]_L)/2^d$. The operators $[\cdot]_L$ and $[\cdot]_M$ correspond to the *least* or *most* significant bits, respectively. We extend the definitions to polynomials by applying $[\cdot]_L$ and $[\cdot]_M$ to each coefficient.

The secret key sk is a tuple of three polynomials $s, e_1, e_2 \leftarrow_{\sigma} \mathcal{R}$, where the entries of e_1 and e_2 have to be small enough for the scheme to be correct; the public key pk consists of the polynomials $a_1, a_2 \leftarrow_{\S} \mathcal{R}_q, b_1 = a_1 s + e_1 \pmod{q}$, and $b_2 = a_2 s + e_2 \pmod{q}$. We depict the algorithm to generate a signature (c', z) for message μ in Figure 1. For verification of the signature (c', z) , it is checked that $z \in \mathcal{R}_{q,[B-U]}$ and that c' equals $H([a_1 z - b_1 c]_M, [a_2 z - b_2 c]_M, \mu)$, with $c = F(c')$. The parameters proposed for ring-TESLA [3] are currently not supported by a security reduction, since in November 2016 an error was detected in the existing reduction. However, we expect that existing security reductions for ring-TESLA's predecessor TESLA [4, 6] are applicable to ring-TESLA as

Sign($\mu; a_1, a_2, s, e_1, e_2$) :

1 $y \leftarrow_{\$} \mathcal{R}_{q,[B]}$ 2 $v_1 \leftarrow a_1 y \pmod{q}$ 3 $v_2 \leftarrow a_2 y \pmod{q}$ 4 $c' \leftarrow H([v_1]_M, [v_2]_M, \mu)$ 5 $c \leftarrow F(c')$ 6 $z \leftarrow y + sc$	7 $w_1 \leftarrow v_1 - e_1 c \pmod{q}$ 8 $w_2 \leftarrow v_2 - e_2 c \pmod{q}$ 9 If $[w_1]_L, [w_2]_L \notin \mathcal{R}_{q,[2^d-L]} \vee z \notin \mathcal{R}_{q,[B-U]}$: 10 Restart 11 Return (z, c')
--	---

Fig. 1. Specification of the scheme ring-TESLA [3].

well². Our modifications described in Section 3.1 and 4 do not depend on the *values* of the parameters. Hence, our results can be applied to other parameter sets of ring-TESLA as long as all values can be represented by the data type `int`.

2.3 Cache side channels

A cache is a small piece of memory that stores selected entries from main memory for quick access by the Central Processing Unit (CPU). If the CPU accesses a memory entry, the access can lead to a cache hit (if the entry is stored in the cache) or to a cache miss (if the entry is not stored in the cache). Inside the cache, the memory entries are stored in sections called cache lines. The sequence of cache lines in a cache is partitioned into cache sets. In a k -way set-associative cache, each cache set consists of k cache lines. A cache has a strategy for replacing entries if the cache is full. A popular strategy is to replace the least recently used entry (LRU strategy). Variants of LRU are used, e.g., in Intel processors [1].

A cache-side-channel vulnerability exists if the interaction between a program and the cache depends on secret information, e.g., on a cryptographic key. In this case, an attacker, observing aspects of this interaction, might learn secret information. Attacks on cryptographic implementations have exploited secret-dependence in the trace of cache hits and misses [2], the time taken for cache hits and misses [8], and the final cache state of an execution [26].

2.4 Leakage bounds on cache side channels

In this article, we follow an approach based on information theory and reachability analysis to compute upper bounds on the cache-side-channel leakage of ring-TESLA (compiled to an x86 binary). Let Obs^a be the set of possible observations an attacker a can make about a single run of an x86 binary. Then $\log_2 |Obs^a|$ is an upper bound on the leakage of the binary with respect to min-entropy [32] and Shannon entropy [30] by [19, Theorem 1] and [5, Theorem 5.3].

We consider the following attacker models $a \in \{acc, accd, trace, time\}$ [13]:

² Security properties of schemes over standard lattices (like TESLA) often hold for corresponding schemes over ideal lattices (like ring-TESLA), e.g., the security reduction from [15] holds for the standard-lattice variant and for the ideal-lattice variant.

accd generalizes techniques like EVICT+TIME and PRIME+PROBE [26]. More concretely, it captures attackers who can determine the number of memory blocks in each cache set in the final cache state after a program execution. *acc* captures attackers who can determine the position of each memory block in the final cache state, inspired by techniques like FLUSH+RELOAD [33]. *trace* captures trace-based attackers who can determine the trace of cache hits and misses that occur during one program execution. For instance, the trace-based attack in [2] uses such traces of hits and misses. *time* models time-based attackers who can observe the running time of one program execution. Actual running times, as used in attacks like [8], are modeled by the amount of cache hits and cache misses that occur.

The possible observations under an attacker model can be computed by reachability analysis [13]. Let \mathcal{D} be the set of the possible states during an execution and let $\text{upd}_{\mathcal{D}} : \mathcal{D} \rightarrow \mathcal{D}$ model the concrete semantics of x86 instructions, i.e., how the execution of instructions updates the state. The possible attacker observations depend on the states that an execution can reach according to $\text{upd}_{\mathcal{D}}$.

Instead of implementing a reachability analysis from scratch, we extend the existing tool CacheAudit 0.2b [21] – a version of CacheAudit [13]. CacheAudit performs a reachability analysis using abstract interpretation [10].

For an abstract reachability analysis, an abstract domain $\overline{\mathcal{D}}$ is defined, which abstracts from details of the concrete execution that are not relevant for the analysis. An abstraction function and a concretization function are defined to convert states between \mathcal{D} and $\overline{\mathcal{D}}$. To represent executions in the abstract domain $\overline{\mathcal{D}}$, an abstract semantics $\text{upd}_{\overline{\mathcal{D}}} : \overline{\mathcal{D}} \rightarrow \overline{\mathcal{D}}$ is defined. To allow a transfer of analysis results from the abstract domain to the concrete domain, the abstract semantics $\text{upd}_{\overline{\mathcal{D}}}$ should be sound with respect to the concrete semantics $\text{upd}_{\mathcal{D}}$, i.e., it should overapproximate the set of reachable states in an execution.

CacheAudit uses multiple abstract domains [13]. The position of a memory block in the cache is abstracted by a set of possible positions. The values of registers and memory entries are abstracted by sets of possible values.

3 Enabling the automatic analysis of ring-TESLA

3.1 Integer implementation of ring-TESLA

We obtained the implementation of ring-TESLA [3] from the authors. The original implementation makes use of floating point operations. CacheAudit 0.2b, however, cannot analyze floating point operations and can therefore not be used to analyze the original implementation directly. An extension of CacheAudit 0.2b to support floating point instructions is out of scope for this article. Changing from floating point to integer operations does not affect the security of the signature scheme, since all operations during the signature generation (cf. Figure 1) are over \mathbb{Z}_q . Moreover, it is a step towards a ring-TESLA implementation for devices without floating point unit, e.g., embedded devices. We replaced all

Listing 3.1. Signature generation in `crypto_sign`

```
1  [...]
2  while(1) {
3      sample_y(vec_y);
4      poly_mul_fixed(vec_v1, vec_y, poly_a1);
5      poly_mul_fixed(vec_v2, vec_y, poly_a2);
6      random_oracle(c, vec_v1, vec_v2, m, mlen);
7      generate_c(pos_list, c);
8
9      computeEc(E1c, sk+sizeof(int)*PARAM_N, pos_list);
10     poly_sub(vec_v1,vec_v1, E1c);
11     if (test_w(vec_v1) != 0){ continue; }
12
13     computeEc(E2c, sk+sizeof(int)*PARAM_N*2, pos_list);
14     poly_sub(vec_v2,vec_v2, E2c);
15     if (test_w(vec_v2) != 0){ continue; }
16
17     computeEc(Sc, sk, pos_list);
18     poly_add(vec_y, vec_y, Sc);
19     if (test_rejection(vec_y) != 0){ continue; }
20
21     for(i=0; i<mlen; i++){ sm[i]=m[i]; }
22     *smLen = CRYPTO_BYTES + mlen;
23     compress_sig(sm+mlen, c, vec_y);
24     return 0; }
```

floating point instructions by integer instructions. The resulting implementation can be analyzed with CacheAudit 0.2c (our extension of CacheAudit 0.2b).

Listing 3.1 shows the parts of the signature generation function `crypto_sign` that are most important for our analysis, leaving out variable declarations.

3.2 Extension of CacheAudit 0.2b

The implementation of the scheme ring-TESLA is the first implementation of post-quantum cryptography (and of lattice-based cryptography) that is analyzed with CacheAudit. The implementation of ring-TESLA contains x86 instructions that are not supported by CacheAudit 0.2b. We extended CacheAudit 0.2b to CacheAudit 0.2c by adding support for these instructions.

To add support for additional x86 instructions to CacheAudit 0.2b, the underlying abstract semantics $\text{upd}_{\overline{\mathcal{D}}}$ must be extended. We implemented abstract semantics for the instructions in the ring-TESLA binary that are unsupported in CacheAudit 0.2b. Table 1 lists the opcodes (unique identifiers) and mnemonics (human-readable descriptions) of the instructions that we added.³

We illustrate the process of extending $\text{upd}_{\overline{\mathcal{D}}}$ at the example of the instruction `Bsr` (Bit scan reverse), which takes the operands `dst` and `src`. The concrete semantics $\text{upd}_{\mathcal{D}}$ of `Bsr dst src` is to compute the index of the most significant bit that is set, i.e., non-zero, in `src` [11]. If such a bit exists in `src`, its index is written to `dst` and the zero flag is set to 0. Otherwise, the zero flag is set to 1.

³ The instructions 0xF7/3 and 0x99 were integrated independently but concurrently into a different version of CacheAudit by Doychev [14].

Type	Opcodes (and mnemonics) of additional instructions
Arithmetic	13 (Adc), 1B (Sbb), 6B (Imul), F7/3 (Neg), F7/4 (Mul), F7/5 (Imul)
Bit string	0FBD (Bsr), 99 (Cdq)
Move	0F9C (Setl), 0F9F (Setg)

Table 1. Additional instructions for ring-TESLA in CacheAudit 0.2c

To support **Bsr**, we extended the parser, the internal instruction representation, and the abstract semantics in CacheAudit 0.2b. We extended the parser to create a **Bsr** instruction in the internal representation when it encounters the opcode 0FBD. We implemented the abstract semantics of **Bsr** by a function `bsr` in the module `valAD`. The function consists of roughly 100 lines of OCaml code.

The function `bsr` operates on sets of potential values for `dst` and `src` and returns a map from possible resulting status flag combinations to the resulting values of registers and memory entries, for which the flag combinations can occur. For each possible value of `src`, we proceed according to the formalization of **Bsr** by Degenbaev [11]. We check whether the value consists only of zeros. In this case, we add a binding (mapping a flag combination to register and memory values) to the resulting map, in which the zero flag is 1 and the value of `dst` is unchanged. Otherwise, we first compute the number of leading zeros by divide and conquer, where we check recursively whether the first half of each non-zero prefix contains bits that are set to 1. The index of the most significant set bit is 64 minus the number of leading zeros. In this case, we add a binding to the resulting map, in which the zero flag is 0 and `dst` contains the computed index.

Our implementation for the other instructions follows the same pattern of parsing and abstract semantics, reusing existing support for similar instructions (e.g., with the same mnemonic) in CacheAudit 0.2b when possible.

4 Detection of potential leakage

We use CacheAudit 0.2c to analyze the signature generation in ring-TESLA for potential leakage of the secret key. We assume that the random number generator is secure and analyze the remaining computation with a few adaptations that allow a meaningful analysis with CacheAudit 0.2c. In the following, we provide details on the configurations of CacheAudit and ring-TESLA, details of our adaptations, and the results of our analysis.

Configuration of CacheAudit. We configure CacheAudit to use a 32kByte, 8-way set-associative data cache with a cache line size of 64Byte. This cache configuration is, e.g., used in the first level cache of the Intel Skylake architecture [18]. As the replacement strategy, we fix LRU.⁴

Configuration of ring-TESLA. We set the parameters of the ring-TESLA scheme to $\text{PARAM_N} = n = 512$, $\text{PARAM_SIGMA} = \sigma = 48$, $\text{PARAM_Q} = q = 33550337$,

⁴ We also investigated FIFO (first in first out) replacement. The leakage bounds (on the unmitigated implementation) are less than 10bit lower than under LRU.

Listing 4.1. Code of the subroutine `generate_c`

```
1 void generate_c(uint32_t *pos_list, unsigned char *c_bin){
2   int32_t c[PARAM_N]; int cnt =0; int pos; [...]
3   crypto_stream(r, R_LENGTH, nonce, c_bin);
4
5   for(i=0; i<PARAM_N; i++){ c[i] = 0;}
6   i=0;
7   while(i<PARAM_W){
8     pos = 0;
9     pos = (r[cnt]<<8) | (r[cnt+1]);
10    pos &= PARAM_N-1;
11    cnt += 2;
12    if (c[pos] == 0) { pos_list[i] = pos; c[pos]=1; i++; cnt++; } }
```

Listing 4.2. Implementation of `sample_y`

```
1 //original
2 do {[...] if(val<0x7ffff) mat_y[i++] = val-PARAM_B; [...]} while(i<PARAM_N);
3
4 //adapted
5 for (i = 0; i < PARAM_N; ++i) { mat_y[i] = *(int *) (0x4) - PARAM_B; }
```

$\text{PARAM_B} = B = 4194303$, $\text{PARAM_W} = \omega = 19$, $\text{PARAM_D} = d = 23$, and $\text{PARAM_U} = U = 2848$. We analyze the function `crypto_sign` from the file `sign.c` in a 32-bit x86 binary of the ring-*TESLA* implementation. To this end, we use a wrapper function that calls `crypto_sign` with an uninitialized secret key, uninitialized message, uninitialized signature buffer, the message size 59 (as in the ring-*TESLA* test suite), and a pointer to `smlen` to store the length of the signed message including the signature. By leaving the secret key and message uninitialized, we treat them as secret input in our analysis.⁵

We compiled the ring-*TESLA* sources and our wrapper with *gcc* version 4.8.4, using `-static` for static linking, `-m32` to target an Intel i386 CPU architecture, and `-fno-stack-protector` to avoid insertion of code for overflow protection.

*Adaptation of ring-*TESLA*.* CacheAudit 0.2c does not support memory accesses that could refer to any possible address, e.g., in the analysis of loop counters that are advanced only under certain conditions and used to index array accesses. This occurs in the ring-*TESLA* routines `generate_c` and `sample_y`.

Listing 4.1 shows the implementation of the function `generate_c`. It uses rejection sampling to generate random values for the parameter array `pos_list`. The loop counter is only increased if the generated value is not rejected. To allow a meaningful analysis, we remove the check that rejects if the same value would occur twice in `pos_list` (highlighted in gray). That is, we overapproximate the possible values of `pos_list`. This manual overapproximation of the semantics preserves the validity of analysis results because it cannot decrease the number of possible attacker observations.

⁵ Treating the message as secret is overly conservative in a signature scenario. We investigated the effect of fixing the message to all '0's and obtained the same leakage bounds as for an uninitialized message in the unmitigated implementation.

Attacker model	<i>acc</i>	<i>accd</i>	<i>trace</i>	<i>time</i>
Leakage in bit	12.9	2.6	51.6	9.5

Table 2. Upper bounds on the leakage of the signature generation

Listing 4.2 shows our adaptation of a loop in `sample_y`. Again, the loop counter advances only if the random number generated in the current iteration satisfies certain criteria. To allow a meaningful analysis with CacheAudit 0.2c, we remove the check of the random number and assign an uninitialized value that overapproximates the possible range to each entry in `mat_y`. With this adaptation, our analysis uses a safe overapproximation of the values that `sample_y` can return, but assumes that `sample_y` itself, i.e., the random number generator, does not have any cache-side-channel leakage.

The function `crypto_sign` contains a potentially infinite while loop, on which CacheAudit 0.2b does not terminate within reasonable time. To make the analysis of this loop feasible, we fix the number of iterations while keeping the source code in the loop body unchanged. More concretely, we fixed the number of iterations to two to account for the effect of more than one iteration.

Note that we use the modifications described in this section only for the initial automatic analysis of ring-TESLA. In the detailed manual inspection in Section 5, we use the unmodified integer implementation.

Analysis Results. We obtain the leakage bounds listed in Table 2.⁶ The bounds lie between 2.6bit and 51.6bit for the different attacker models. One run of the adapted ring-TESLA leaks at most 2.6bit to attackers under *accd*, at most 9.5bit to attackers under *time*, at most 12.9bit to attackers under *acc*, and at most 51.6bit to attackers under *trace*. In the remainder of this article we investigate whether these non-zero leakage bounds are substantiated by concrete threats and how the leakage bounds can be reduced by mitigating concrete threats.

5 Manual analysis of the potential leakage

We manually analyze the signature generation `crypto_sign` to check if the potential leakage detected by program analysis corresponds to an actual concern. We identify substantiated threats of leakage to cache side channels (CSCs) in the routines `generate_c`, `test_w`, `test_rejection`, and `computeEc`.⁷ The following variables have to be kept secret during the execution of `crypto_sign`: the secret key `sk`, the randomness `vec_y`, the polynomials `vec_v1` and `vec_v2` to compute the hash value, and the polynomials `E1c`, `E2c`, and `Sc`. Furthermore, the hash value `c` and the representation of the corresponding encoded polynomial `pos_list` have to be kept secret until line 19 in Listing 3.1. In line 19, it

⁶ Throughout the article, we round bounds up to one decimal place and truncate the bounds to $3 * \text{PARAM_N} * 32\text{bit} = 49152\text{bit}$, i.e., the maximum size of the key.

⁷ The analysis of the other subroutines of `crypto_sign` can be found in the corresponding technical report under <http://eprint.iacr.org/2017/951>.

Listing 5.1. Code of the subroutine `computeEc`

```
1 static void computeEc(poly Ec, const unsigned char *sk, const
2 uint32_t pos_list[PARAM_W]) {
3     int i, j, pos, * e;
4     e = (int*)sk;
5     for(i=0; i<PARAM_N; i++){ Ec[i] = 0;}
6
7     for(i=0; i<PARAM_W; i++){
8         pos = pos_list[i];
9         for(j=0; j<pos; j++){ Ec[j] += e[j+PARAM_N - pos];}
10        for(j=pos; j<PARAM_N; j++){ Ec[j] -= e[j-pos];} } }
```

Listing 5.2. Code of the subroutine `test_rejection`

```
1 static int test_rejection(poly poly_z) {
2     int i;
3     for(i=0; i<PARAM_N; i++){
4         if(poly_z[i]<-(PARAM_B-PARAM_U) || poly_z[i]>(PARAM_B-PARAM_U)){return 1;}}
5     return 0; }
```

is decided whether the potential signature (computed in line 18) and `c` are returned and, hence, whether `c` and `pos_list` become public information (via the encoding function F , the values can be computed from each other), or whether all computed values are discarded. An attacker should not learn the values of the discarded polynomials, e.g., `c` or `pos_list`. If the attacker learns values of `pos_list` or `c`, there exists a potential attack as described below.

Analysis of the subroutine `compute_Ec`. The implementation of the subroutine `compute_Ec` is given in Listing 5.1. The values that have to be kept secret during this computation are `sk`, `e`, `pos_list`, and `pos`. Most loops and branchings do not depend on any of the secret values. However, there might be a possible leakage of `pos` (and hence of the values in the secret `pos_list`) because of the cache hits/misses depending on `e`. In both loop bodies values are read from `e` (namely, either `e[j+PARAM_N-pos]` or `e[j-pos]`) such that in both loops together all entries of `e` are read. However, leakage arises from the chronological order of cache hits and misses. We illustrate the leakage using an example: The array `e` consists of `PARAM_N` many entries of type `int`, i.e., each entry of `e` is represented in 32bit. Since one cache line is 64Byte (cf. Section 4), 16 entries (depending on the alignment in the memory) of `e` fit into one cache line. Let `pos=14`. Then, under the trace-driven attacker model, an attacker sees one cache miss (on element `e[PARAM_N-14]`) and 13 cache hits (on elements `e[PARAM_N-13]`, ..., `e[PARAM_N-1]`) during the loop in line 9.⁸ However, two more entries of `e` are also already loaded in the cache, namely `e[PARAM_N-16]` and `e[PARAM_N-15]`. Thus, in the second loop in line 10, the attacker sees cache hits on those two elements. He might, hence, be able to determine the value of `pos` from the distribution of the hits for the considered cache line over the loops.

⁸ To simplify our explanation we assume that the corresponding cache line starts with `e[PARAM_N-16]` and ends with `e[PARAM_N-1]`.

Analysis of the subroutine `test_rejection`. The implementation of the subroutine `test_rejection` is given in Listing 5.2. The variable `poly_z` in Listing 5.2 has to be kept secret. The subroutine `test_rejection` consists of a for-loop that loops independently of the secret over $i=0, \dots, \text{PARAM_N}$. Within the for-loop, there is a secret-dependent if-condition. This leads to a CSC vulnerability.

Assume a strong trace-driven attacker model, i.e., the attacker has a sequence of occurred cache hits and misses. Assume furthermore that `poly_z` is already loaded in the cache before the if-condition is evaluated.⁹ When the if-condition in line 4, Listing 5.2, is never true, then the value 0 is returned and the attacker gets a sequence of `PARAM_N` (or $2 \cdot \text{PARAM_N}$ — depending on the compilation) hits. This essentially means that all coefficients of `poly_z` are in the interval $[-B + U, B - U]$ and, hence, the corresponding signature is compressed and returned (cf. Listing 3.1). Next, we consider the other case, i.e., the absolute value of at least one of the coefficients of `poly_z` is larger than $B - U$. That means that the if-condition in Listing 5.2 holds true for some $i \in \{0, \dots, \text{PARAM_N}\}$. Hence, 1 is returned in the i -th iteration and the attacker gets a sequence of only `PARAM_N - i` hits. Hence, the attacker knows the exact index of the coefficient that violated the if-condition. Assume the attacker also knows the values in the array `pos_list` (which corresponds to the polynomial $c = F(H([v_1]_M, [v_2]_M, \mu))$ in Figure 1) from another cache-side-channel vulnerability. Then the attacker might know which coefficients of the secret s contributed to the i -th, large coefficient of `poly_z`. If an attacker learns the exact position i and the corresponding `pos_list` for many different values to the same secret key s then the attacker might receive enough information about the size of the entries in s to successfully break the scheme via a learning-the-parallelepiped-attack [16, 24].

Analysis of the subroutine `test_w`. The implementation of the routine `test_w` is given in Listing 5.3. The values `poly_w`, `val`, and `left` in Listing 5.3 have to be kept secret. The CSC vulnerability is similar to the channel described previously for `test_rejection`. In the subroutine `test_w`, the CSC vulnerability comes from the early abortion depending on `left` in line 12 of Listing 5.3. When the if-condition in line 12 holds for some i and the corresponding `abs(left)`, -1 is immediately returned and a trace-driven attacker might learn the exact index i .

Analysis of the subroutine `generate_c`. The implementation of the subroutine `generate_c` is given in Listing 4.1. The values `pos_list`, `c`, and `pos` in Listing 4.1 have to be kept secret. There are no branchings or loops depending on the secret value, except for one if-condition on `c[pos]` in line 12 of Listing 4.1. If a cache with no-write-allocate policy is used, the values `c[i]` are not cached in line 5. Hence, an attacker might be able to find out which elements `c[i]` are cached in line 12 and to learn information about the values of `pos`. Together with the vulnerability in `test_rejection`, an attacker might be able to successfully break the scheme via a learning-the-parallelepiped-attack [16, 24].

⁹ Our arguments hold also true if we assume that `poly_z` is not loaded in the cache. In the ring-TESLA implementation, it is already loaded in line 18 in Listing 3.1.

Listing 5.3. Code of the subroutine `test_w`

```
1 static int test_w(poly poly_w)
2 { int i; int64_t left, right, val;
3   for(i=0; i<PARAM_N; i++){
4     val = (int64_t) poly_w[i];
5     val = val % PARAM_Q;
6     if (val < 0){ val = val + PARAM_Q;}
7     left = val;
8     left = left % (1<<(PARAM_D));
9     left -= (1<<PARAM_D)/2;
10    left++;
11    right = (1<<(PARAM_D-1))-PARAM_REJECTION;
12    if (abs(left) > right){ return -1; }
13    return 0; }
```

Combined analysis of the overall signature generation. The most important parts of the implementation of `crypto_sign` are depicted in Listing 3.1. In the signature generation, most operations, branchings, or loops are independent of the *value* of the secret. Exceptions are the branchings in lines 11, 15, and 19 in Listing 3.1: They depend on secret values and, hence, the length of the observed trace of cache hits and misses depends on the branches that are taken. What does this mean from a cryptographic viewpoint? Assuming the subroutine `test_w` does not leak any bit, then the attacker does not learn more information about the secret if he knows whether or not the condition in line 11 holds. The attacker would just learn that `vec_v1` does not fulfill the conditions needed for a valid signature. However, the attacker does not learn why exactly the condition was not fulfilled (the attacker does not learn the index on which the if-condition failed). Furthermore, since the value `vec_v1` depends on `vec_y` and the value `vec_y` is discarded if the if-condition in line 11 does not hold, the attacker does not get any additional information about the secret he did not know before. The same explanation also holds for the branchings in line 15 and line 19.

In summary, this means that there exists a potential leakage that we probably cannot get rid of, but it does not affect the security of the signature scheme as long as we do not have leakages in `test_w` and `test_rejection` or `generate_c`.

6 Mitigation of the vulnerabilities

6.1 Adaptation of vulnerable routines

In Section 5, we identified substantiated threats of CSC leakage in the routines `test_w`, `test_rejection`, `computeEc`, and `generate_c`. Since the leakage in `generate_c` is only a concern in combination with the leakage in `test_w` and `test_rejection`, it suffices to analyze and mitigate the leakage in `test_w`, `test_rejection`, and `computeEc`. We analyze `test_w`, `test_rejection`, and `computeEc` individually with CacheAudit 0.2c to obtain leakage bounds on the unmitigated implementations. The leakage bounds are listed in Table 3. There are, indeed, non-zero bounds for all three unmitigated routines.

	Unmitigated routines				Mitigated routines			
	<i>acc</i>	<i>accd</i>	<i>trace</i>	<i>time</i>	<i>acc</i>	<i>accd</i>	<i>trace</i>	<i>time</i>
<code>test_w</code>	31	31	49152	19.3	0	0	0	0
<code>test_rejection</code>	31	31	10.1	10.1	0	0	0	0
<code>computeEc</code>	0	0	20	5.9	0	0	19	4.4
<code>crypto_sign</code>	12.9	2.6	51.6	9.5	8.1	1.6	48.6	9.0

Table 3. Leakage bounds [bit]

```

1 int test_rejection(poly poly_z) {
2   int i; int res; res = 0;
3   for(i=0; i<PARAM_N; i++){
4     res |= (poly_z[i] < -(PARAM_B-PARAM_U));
5     res |= (poly_z[i] > (PARAM_B-PARAM_U)); }
6   return res; }
7
8 int test_w(poly poly_w) { [...]
9   for(i=0; i<PARAM_N; i++) {
10    val = poly_w[i]; val = val % PARAM_Q;
11    val += (((unsigned int)val & 0x80000000) > 31)*PARAM_Q;
12    left = val; left = left % (1<<(PARAM_D));
13    left -= (1<<(PARAM_D)/2); left++;
14    right = (1<<(PARAM_D-1))-PARAM_REJECTION;
15    res |= (abs(left) - right > 0); }
16   return -res; }
17
18 void computeEc([...]) { [...]
19   for(i=0; i<PARAM_N; i++) Ec[i] = 0;
20   for(i=0; i<PARAM_N; i++) tmp = e[i];
21   for(i=0; i<PARAM_W; i++) {
22     pos = pos_list[i];
23     for(j=0; j<pos; j++) { Ec[j] += e[j+PARAM_N - pos]; }
24     for(j=pos; j<PARAM_N; j++) { Ec[j] -= e[j-pos]; } } }

```

We adapt the routines, as shown in the above listing, to mitigate the leakage. In `test_rejection`, we collect the result, i.e., whether 0 or 1 is returned, in an auxiliary variable `res` and return it after `PARAM_N` iterations, instead of returning early in case of a failed test. In `test_w`, we also collect the result, i.e., whether 0 or 1 is returned, in an auxiliary variable `res`, instead of returning early in case of failure. Furthermore, we replace the branching on `val` by an assignment that masks the value by the branching condition. The idea to mask assignments by branching conditions comes from conditional assignment [22] - a program transformation to mitigate timing side channels, which performs rather well in practical evaluation [20]. In `computeEc`, we add preloading of the variable `e` to ensure that the sequence of cache hits and misses does not depend on the secret-dependent order of accesses (under the assumption that no process interferes with the cache during the ring-TESLA execution).

By code inspection, the modifications should remove the CSC leakage in the three routines. In the following, we investigate this with CacheAudit.

6.2 Analysis of the effectiveness of the mitigations

We analyze the mitigated routines `test_w`, `test_rejection`, and `computeEc` using CacheAudit 0.2c and obtain the leakage bounds listed in Table 3. For `test_w` and `test_rejection`, we obtain the leakage bound 0bit for all attacker models. Thus, we effectively removed the potential leakage. For `computeEc`, we obtain 0bit leakage bounds for `acc` and `accd` and non-zero leakage bounds, namely 19bit and 4.4bit, for `trace` and `time`, respectively. The bounds computed with CacheAudit are provable upper bounds, but not necessarily tight. The preloading of `e` should remove the leakage from `computeEc`, because it makes the caching of `e` independent of secrets. Since CacheAudit was able to recognize preloading as effective in other cases [13, 21], it is interesting why CacheAudit 0.2c does not yield a 0bit leakage bound in this case. The investigation and fine-tuning of the analysis precision is an interesting direction for future work.

Table 3 lists also the leakage bounds we obtain on `crypto_sign` with and without our countermeasures. All four leakage bounds are reduced by our countermeasures. The highest reduction is achieved for the `acc` leakage bound. The `acc` leakage bound is reduced by 4.8bit to 8.1bit.

Based on our manual inspection of the individual routines in ring-TESLA, there are two possible sources for the remaining potential leakage reported by CacheAudit 0.2c. One source is `generate_c`, where, as discussed above, the remaining leakage is harmless. The second source is the rejection sampling in `crypto_sign`. This matches the fact that the leakage bounds are non-zero. Note that, since we compute upper bounds on the leakage based on overapproximation, the actual leakage of the implementation could be even lower than the reported leakage bounds. We expect in particular the leakage bounds for `trace` and `time` to be quite conservative because CacheAudit 0.2c was not able to recognize the preloading countermeasure in the implementation of `computeEc` and the results of `computeEc` are propagated further through the implementation. Nevertheless, the bounds show that the CSC leakage of ring-TESLA to `acc`, `accd`, `trace`, and `time` is rather low. For `acc`, `accd`, and `time`, it even lies below 10bit.

Note that, the leakage bounds refer to information about the secret key, i.e., about the three polynomials that all together are saved in roughly 38,400bit. By construction of the learning with errors problem (LWE) — the underlying hardness assumption of ring-TESLA — the potential leakage of at most 49bit of the secret key does not immediately translate to the bit-hardness of LWE (resp., the bit-security of ring-TESLA).

7 Conclusion

In this article, we analyzed an implementation of the lattice-based signature scheme ring-TESLA for cache-side-channel vulnerabilities. We identified four routines in the implementation that are vulnerable through cache side channels. Two of these routines, a rejection sampling and a signature validity check, use a secret-dependent number of iterations. One routine is a sparse polynomial multiplication that traverses one polynomial in a secret-dependent order. We

modified these functions to ensure a constant number of iterations and secret-independent caching of the polynomial. By modifying these functions, we also eliminated the possibility to exploit the fourth vulnerability. For the modified ring-TESLA implementation, we obtained low upper bounds on the leakage to four attacker models, using program analysis.

Our results show that implementations of rejection sampling and sparse multiplication should be inspected for side channels with particular care. While these techniques are not very common in classical cryptography like RSA, they play a significant role in post-quantum cryptography. Rejection sampling occurs, for instance, in multiple lattice-based signature schemes [6, 7, 15] and in key exchange protocols [35]. Sparse multiplication also occurs in many lattice-based schemes [7, 9, 15]. Overall, the implementation and analysis of post-quantum cryptography poses additional challenges compared to classical cryptography.

Acknowledgements We thank the anonymous reviewers for their helpful suggestions and Sedat Akleylek for contributing to our modifications of the original ring-TESLA implementation. This work has been partially funded by the DFG as part of projects P1 and E3 within the CRC 1119 CROSSING.

References

1. Abel, A., Reineke, J.: Reverse engineering of cache replacement policies in Intel microprocessors and their evaluation. In: ISPASS. pp. 141–142 (2014)
2. Aciçmez, O., Ç. K. Koç: Trace-driven cache attacks on AES. Cryptology ePrint Archive, Report 2006/138
3. Akleylek, S., Bindel, N., Buchmann, J., Krämer, J., Marson, G.A.: An Efficient Lattice-Based Signature Scheme with Provably Secure Instantiation. In: AFRICACRYPT. pp. 44–60 (2016)
4. Alkim, E., Bindel, N., Buchmann, J., Dagdelen, Ö., Eaton, E., Gutoski, G., Krämer, J., Pawlega, F.: Revisiting TESLA in the quantum random oracle model. In: PQCrypto. pp. 143–162 (2017)
5. Alvim, M.S., Chatzikokolakis, K., Palamidessi, C., Smith, G.: Measuring Information Leakage using Generalized Gain Functions. In: CSF. pp. 265–279 (2012)
6. Bai, S., Galbraith, S.D.: An Improved Compression Technique for Signatures Based on Learning with Errors. In: CT-RSA. pp. 28–47 (2014)
7. Barreto, P.S.L.M., Longa, P., Naehrig, M., Ricardini, J.E., Zanon, G.: Sharper Ring-LWE Signatures. Cryptology ePrint Archive, Report 2016/1026
8. Bernstein, D.J.: Cache-timing attacks on AES. Tech. rep., University of Illinois at Chicago (2005)
9. Buchmann, J., Göpfert, F., Güneysu, T., Oder, T., Pöppelmann, T.: High-Performance and Lightweight Lattice-Based Public-Key Encryption. In: IoTPTS. pp. 2–9 (2016)
10. Cousot, P., Cousot, R.: Abstract Interpretation: A Unified Lattice Model for Static Analysis of Programs by Construction or Approximation of Fixpoints. In: POPL. pp. 238–252 (1977)
11. Degenbaev, U.: Formal Specification of the x86 Instruction Set Architecture. Ph.D. thesis, Universität des Saarlandes (2012)
12. Doychev, G., Köpf, B.: Rigorous Analysis of Software Countermeasures against Cache Attacks. In: PLDI. pp. 406–421 (2017)

13. Doychev, G., Köpf, B., Mauborgne, L., Reineke, J.: CacheAudit: A Tool for the Static Analysis of Cache Side Channels. *ACM TISSEC* 18(1), 4:1–4:32 (2015)
14. Doychev, G.: Commit f063813faa548da9bfb11dea9ff6fe39c0f11626: Adding support for CDQ and NEG instructions. <https://github.com/cacheaudit/cacheaudit/commit/f063813faa548da9bfb11dea9ff6fe39c0f11626> (2016), [Online; accessed 05/23/2017]
15. Ducas, L., Durmus, A., Lepoint, T., Lyubashevsky, V.: Lattice Signatures and Bimodal Gaussians. In: *CRYPTO*. pp. 40–56 (2013)
16. Ducas, L., Nguyen, P.Q.: Learning a Zonotope and More: Cryptanalysis of NTRUSign Countermeasures. In: *ASIACRYPT*. pp. 433–450 (2012)
17. Groot Bruinderink, L., Hülsing, A., Lange, T., Yarom, Y.: Flush, Gauss, and Reload – A Cache Attack on the BLISS Lattice-Based Signature Scheme. In: *CHES*. pp. 323–345 (2016)
18. Intel Corporation: Intel[®] 64 and IA-32 Architectures Optimization Reference Manual. Order Number: 248966-032 (2016)
19. Köpf, B., Smith, G.: Vulnerability Bounds and Leakage Resilience of Blinded Cryptography under Timing Attacks. In: *CSF*. pp. 44–56 (2010)
20. Mantel, H., Starostin, A.: Transforming Out Timing Leaks, More or Less. In: *ESORICS*. pp. 447–467 (2015)
21. Mantel, H., Weber, A., Köpf, B.: A Systematic Study of Cache Side Channels across AES Implementations. In: *ESSoS*. pp. 213–230 (2017)
22. Molnar, D., Piotrowski, M., Schultz, D., Wagner, D.: The Program Counter Security Model: Automatic Detection and Removal of Control-Flow Side Channel Attacks. In: *ICISC*. pp. 156–168 (2005)
23. National Institute of Standards and Technology (NIST): Post-quantum project. https://pqcrypto2016.jp/data/pqc2016_nist_announcement.pdf (2016), [Online; accessed 05/23/2017]
24. Nguyen, P.Q., Regev, O.: Learning a Parallelepiped: Cryptanalysis of GGH and NTRU Signatures. In: *EUROCRYPT*. pp. 271–288 (2006)
25. Oder, T., Schneider, T., Pöppelmann, T., Güneysu, T.: Practical CCA2-Secure and Masked Ring-LWE Implementation. *Cryptology ePrint Archive*, Report 2016/1109
26. Osvik, D.A., Shamir, A., Tromer, E.: Cache Attacks and Countermeasures: The Case of AES. In: *CT-RSA*. pp. 1–20 (2006)
27. Pessl, P.: Analyzing the Shuffling Side-Channel Countermeasure for Lattice-Based Signatures. *Cryptology ePrint Archive*, Report 2017/033
28. Reparaz, O., Roy, S.S., Vercauteren, F., Verbauwhede, I.: A masked ring-LWE implementation. *Cryptology ePrint Archive*, Report 2015/724
29. Saarinen, M.J.O.: Arithmetic coding and blinding countermeasures for lattice signatures. *Cryptology ePrint Archive*, Report 2016/276
30. Shannon, C.E.: A Mathematical Theory of Communication. *Bell System Technical Journal* 27(3), 379–423, 623–656 (1948)
31. Shor, P.W.: Polynomial-Time Algorithms for Prime Factorization and Discrete Logarithms on a Quantum Computer. *SIAM J. Comput.* 26(5), 1484–1509 (1997)
32. Smith, G.: On the Foundations of Quantitative Information Flow. In: *FOSSACS*. pp. 288–302 (2009)
33. Yarom, Y., Falkner, K.: FLUSH+RELOAD: A High Resolution, Low Noise, L3 Cache Side-Channel Attack. In: *USENIX Security*. pp. 719–732 (2014)
34. Yarom, Y., Genkin, D., Heninger, N.: CacheBleed: A Timing Attack on OpenSSL Constant Time RSA. In: *CHES*. pp. 346–367 (2016)
35. Zhang, J., Zhang, Z., Ding, J., Snook, M., Dagdelen, Ö.: Authenticated key exchange from ideal lattices. In: *EUROCRYPT*. pp. 719–751 (2015)