

# CRYSTALS-Dilithium: A Lattice-Based Digital Signature Scheme

Léo Ducas<sup>1</sup>, Eike Kiltz<sup>2</sup>, Tancrede Lepoint<sup>3</sup>, Vadim Lyubashevsky<sup>4</sup>,  
Peter Schwabe<sup>5</sup>, Gregor Seiler<sup>6</sup> and Damien Stehlé<sup>7</sup>

<sup>1</sup> CWI, Netherlands

<sup>2</sup> Ruhr Universität Bochum, Germany

<sup>3</sup> SRI International, USA

<sup>4</sup> IBM Research – Zurich, Switzerland

<sup>5</sup> Radboud University, Netherlands

<sup>6</sup> IBM Research – Zurich and ETH Zurich, Switzerland

<sup>7</sup> ENS de Lyon, France

**Abstract.** In this paper, we present the lattice-based signature scheme Dilithium, which is a component of the CRYSTALS (Cryptographic Suite for Algebraic Lattices) suite that was submitted to NIST’s call for post-quantum cryptographic standards. The design of the scheme avoids all uses of discrete Gaussian sampling and is easily implementable in constant-time. For the same security levels, our scheme has a public key that is 2.5X smaller than the previously most efficient lattice-based schemes that did not use Gaussians, while having essentially the same signature size. In addition to the new design, we significantly improve the running time of the main component of many lattice-based constructions – the number theoretic transform. Our AVX2-based implementation results in a speed-up of roughly a factor of 2 over the previously best algorithms that appear in the literature. The techniques for obtaining this speed-up also have applications to other lattice-based schemes.

**Keywords:** Lattice Cryptography · Digital Signatures · Constant-Time Implementation · AVX2

## 1 Introduction

Cryptography based on the hardness of lattice problems is seen as a very promising replacement of traditional cryptography after the eventual coming of quantum computers. In this paper, we present a new digital signature scheme Dilithium, whose security is based on the hardness of finding short vectors in lattices. We also present its complete optimized implementation, as well as a detailed security analysis of the underlying problems upon which it is based. Our scheme was designed with the following criteria in mind:

**Simple to implement securely.** The most compact lattice-based signature schemes [DDLL13, DLP14] crucially require the generation of secret randomness from the discrete Gaussian distribution. Generating such samples in a way that is secure against side-channel attacks is highly non-trivial and can easily lead to insecure implementations, as demonstrated in [BHLY16, EFGT17, PBY17]. While it may be possible that a very careful implementation can prevent such attacks, it is unreasonable to assume that a universally-deployed scheme containing many subtleties will always be expertly implemented. Dilithium therefore only uses uniform sampling, as was originally proposed in [Lyu09, GLP12, BG14]. Furthermore all other operations (such as polynomial multiplication and rounding) are easily implemented in constant time.

<pre> Gen 01 <math>\mathbf{A} \leftarrow R_q^{k \times \ell}</math> 02 <math>(\mathbf{s}_1, \mathbf{s}_2) \leftarrow S_\eta^\ell \times S_\eta^k</math> 03 <math>\mathbf{t} := \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2</math> 04 <b>return</b> <math>(pk = (\mathbf{A}, \mathbf{t}), sk = (\mathbf{A}, \mathbf{t}, \mathbf{s}_1, \mathbf{s}_2))</math>  Sign(<math>sk, M</math>) 05 <math>\mathbf{z} := \perp</math> 06 <b>while</b> <math>\mathbf{z} = \perp</math> <b>do</b> 07   <math>\mathbf{y} \leftarrow S_{\gamma_1 - 1}^\ell</math> 08   <math>\mathbf{w}_1 := \text{HighBits}(\mathbf{A}\mathbf{y}, 2\gamma_2)</math> 09   <math>c \in B_{60} := \text{H}(M \parallel \mathbf{w}_1)</math> 10   <math>\mathbf{z} := \mathbf{y} + c\mathbf{s}_1</math> 11   <b>if</b> <math>\ \mathbf{z}\ _\infty \geq \gamma_1 - \beta</math> <b>or</b> <math>\ \text{LowBits}(\mathbf{A}\mathbf{y} - c\mathbf{s}_2, 2\gamma_2)\ _\infty \geq \gamma_2 - \beta</math>, <b>then</b> <math>\mathbf{z} := \perp</math> 12 <b>return</b> <math>\sigma = (\mathbf{z}, c)</math>  Verify(<math>pk, M, \sigma = (\mathbf{z}, c)</math>) 13 <math>\mathbf{w}'_1 := \text{HighBits}(\mathbf{A}\mathbf{z} - c\mathbf{t}, 2\gamma_2)</math> 14 <b>if return</b> <math>\llbracket \ \mathbf{z}\ _\infty &lt; \gamma_1 - \beta \rrbracket</math> <b>and</b> <math>\llbracket c = \text{H}(M \parallel \mathbf{w}'_1) \rrbracket</math> </pre>
--

**Figure 1:** Template for our signature scheme.

**Be conservative with parameters.** Since we are aiming for long-term security, we have analyzed the applicability of lattice attacks from a very favorable, to the attacker, viewpoint. In particular, we are considering quantum algorithms that require virtually as much space as time. Such algorithms are currently unrealistic, and there seem to be serious obstacles in removing the space requirement, but we are allowing for the possibility that improvements may occur in the future.

**Minimize the size of public key + signature.** Since many applications require the transmission of both the public key and the signature (e.g. certificate chains), we designed our scheme to minimize the sum of these parameters. Under the restriction that we avoid discrete Gaussian sampling, to the best of our knowledge, Dilithium has the smallest combination of signature and public key sizes of any post-quantum signature scheme.

**Be modular – easy to vary security.** The two operations that constitute nearly the entirety of the signing and verification procedures are expansion of an XOF (we use SHAKE-128 and SHAKE-256), and multiplication in the polynomial ring  $\mathbb{Z}_q[X]/(X^n + 1)$ . Highly efficient implementations of our algorithm will therefore need to optimize these operations and make sure that they run in constant time. For all security levels, our scheme uses the same ring with  $q = 2^{23} - 2^{13} + 1$  and  $n = 256$ . Varying security simply involves doing more/less operations over this ring and doing more/less expansion of the XOF. In other words, once an optimized implementation is obtained for some security level, it is almost trivial to obtain an optimized implementation for a higher/lower level.

## 1.1 Overview of the Basic Approach

The design of the scheme is based on the “Fiat-Shamir with Aborts” approach [Lyu09] and bears most resemblance to the schemes proposed in [GLP12, BG14]. For readers who are unfamiliar with the general framework of such signature schemes, we present a simplified (and less efficient) version of our scheme in Fig. 1. This version is essentially a slightly modified version of the scheme from [BG14]. We will now go through each of its components to give the reader an idea of how such schemes work.

**Key Generation.** The key generation algorithm generates a  $k \times \ell$  matrix  $\mathbf{A}$  each of whose entries is a polynomial in the ring  $R_q = \mathbb{Z}_q[X]/(X^n + 1)$ . As previously mentioned, we will always have  $q = 2^{23} - 2^{13} + 1$  and  $n = 256$ . Afterwards, the algorithm samples random secret key vectors  $\mathbf{s}_1$  and  $\mathbf{s}_2$ . Each coefficient of these vectors is an element of  $R_q$  with small coefficients – of size at most  $\eta$ . Finally, the second part of the public key is computed as  $\mathbf{t} = \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2$ . All algebraic operations in this scheme are assumed to be over the polynomial ring  $R_q$ .

**Signing Procedure.** The signing algorithm generates a masking vector of polynomials  $\mathbf{y}$  with coefficients less than  $\gamma_1$ . The parameter  $\gamma_1$  is set strategically – it is large enough that the eventual signature does not reveal the secret key (i.e. the signing algorithm is zero-knowledge), yet small enough so that the signature is not easily forged. The signer then computes  $\mathbf{A}\mathbf{y}$  and sets  $\mathbf{w}_1$  to be the “high-order” bits of the coefficients in this vector. In particular, every coefficient  $w$  in  $\mathbf{A}\mathbf{y}$  can be written in a canonical way as  $w = w_1 \cdot 2\gamma_2 + w_0$  where  $|w_0| \leq \gamma_2$ ;  $\mathbf{w}_1$  is then the vector comprising all the  $w_1$ ’s. The challenge  $c$  is then created as the hash of the message and  $\mathbf{w}_1$ . The output  $c$  is a polynomial in  $R_q$  with exactly  $60 \pm 1$ ’s and the rest 0’s. The reason for this distribution is that  $c$  has small norm and comes from a domain of size  $> 2^{256}$ . The potential signature is then computed as  $\mathbf{z} = \mathbf{y} + c\mathbf{s}_1$ .

If  $\mathbf{z}$  were directly output at this point, then the signature scheme would be insecure due to the fact that the secret key would be leaked. To avoid the dependency of  $\mathbf{z}$  on the secret key, we use rejection sampling. The parameter  $\beta$  is set to be the maximum possible coefficient of  $c\mathbf{s}_1$ . Since  $c$  has  $60 \pm 1$ ’s and the maximum coefficient in  $\mathbf{s}_1$  is  $\eta$ , it’s easy to see that  $\beta \leq 60\eta$ . If any coefficient of  $\mathbf{z}$  is larger than  $\gamma_1 - \beta$ , then we reject and restart the signing procedure. Also, if any coefficient of the low-order bits of  $\mathbf{A}\mathbf{z} - c\mathbf{t}$  is greater than  $\gamma_2 - \beta$ , we restart. The first check is necessary for security, while the second is necessary for both security and correctness. The while loop in the signing procedure keeps being repeated until the preceding two conditions are satisfied. The parameters are set such that the expected number of repetitions is between 4 and 7.

**Verification.** The verifier first computes  $\mathbf{w}'_1$  to be the high-order bits of  $\mathbf{A}\mathbf{z} - c\mathbf{t}$ , and then accepts if all the coefficients of  $\mathbf{z}$  are less than  $\gamma_1 - \beta$  and if  $c$  is the hash of the message and  $\mathbf{w}'_1$ . Let us look at why verification works, in particular as to why  $\text{HighBits}(\mathbf{A}\mathbf{z} - c\mathbf{t}, 2\gamma_2) = \text{HighBits}(\mathbf{A}\mathbf{y}, 2\gamma_2)$ . The first thing to notice is that  $\mathbf{A}\mathbf{z} - c\mathbf{t} = \mathbf{A}\mathbf{y} - c\mathbf{s}_2$ . So all we really need to show is that

$$\text{HighBits}(\mathbf{A}\mathbf{y}, 2\gamma_2) = \text{HighBits}(\mathbf{A}\mathbf{y} - c\mathbf{s}_2, 2\gamma_2). \quad (1)$$

The reason for this is that a valid signature will have  $\|\text{LowBits}(\mathbf{A}\mathbf{y} - c\mathbf{s}_2, 2\gamma_2)\|_\infty < \gamma_2 - \beta$ . And since we know that the coefficients of  $c\mathbf{s}_2$  are smaller than  $\beta$ , we know that adding  $c\mathbf{s}_2$  is not enough to cause any carries by increasing any low-order coefficient to have magnitude at least  $\gamma_2$ . Thus Eq. (1) is true and the signature verifies correctly.

## 1.2 Dilithium

The basic template in Fig. 1 is rather inefficient, as is. The most glaring (but easily fixed) inefficiency is that the public key consists of a matrix of  $k \cdot \ell$  polynomials, which could have a rather large representation. The obvious fix is to have  $\mathbf{A}$  generated from some seed  $\rho$  using SHAKE-128, and this is a standard technique. The novelty of Dilithium over the previous schemes is that we also shrink the size of the public key by a factor of 2.5 at the expense of increasing the signature by around 150 bytes. For the recommended security level, the scheme has 2.7KB signatures and 1.5KB public keys.

The main observation for obtaining this very favorable trade-off is that when the verifier computes  $\mathbf{w}'_1$  in Line 13, the high-order bits of  $\mathbf{Az} - \mathbf{ct}$  do not depend too much on the *low order* bits of  $\mathbf{t}$  because  $\mathbf{t}$  is being multiplied by a very low-weight polynomial  $c$ . In our scheme, some low-order bits of  $\mathbf{t}$  are not included in the public key, and so the verifier cannot always correctly compute the high-order bits of  $\mathbf{Az} - \mathbf{ct}$ . To make up for this, the signer includes some “hints” as part of the signature, which are essentially the carries caused by adding in the product of  $c$  with the missing low-order bits of  $\mathbf{t}$ . With this hint, the verifier is able to correctly compute  $\mathbf{w}'_1$ .

Additionally, we make our scheme deterministic using the standard technique of adding a seed to the secret key and using this seed together with the message to produce the randomness  $\mathbf{y}$  in Line 07. Our full scheme in Fig. 4 also makes use of basic optimizations such as pre-hashing the message  $M$  so as to not rehash it with every signing attempt.

**Implementation Considerations.** The main algebraic operation performed in the scheme is a multiplication of a matrix  $\mathbf{A}$ , whose elements are polynomials in  $\mathbb{Z}_q[X]/(X^{256} + 1)$  by a vector of such polynomials. In our recommended parameter setting,  $\mathbf{A}$  is a  $5 \times 4$  matrix, and thus the multiplication  $\mathbf{Av}$  involves 20 polynomial multiplications. As in most lattice-based schemes that are based on operations over polynomial rings, we have chosen our ring so that the multiplication operation has a very efficient implementation via the Number Theoretic Transform (NTT), which is just a version of FFT that works over the finite field  $\mathbb{Z}_q$  rather than over the complex numbers. To enable the NTT, we needed to choose a prime  $q$  so that the group  $\mathbb{Z}_q^*$  has an element of order  $2n = 512$ ; or equivalently  $q \equiv 1 \pmod{512}$ . If  $r$  is such an element, then  $X^{256} + 1 = (X - r)(X - r^3) \cdots (X - r^{511})$  and thus one can equivalently represent any polynomial  $a \in \mathbb{Z}_q[X]/(X^{256} + 1)$  in its CRT (Chinese Remainder Theorem) form as  $(a(r), a(r^3), \dots, a(r^{2n-1}))$ . The advantage of this representation is that the product of two polynomials is coordinate-wise. Therefore the most expensive parts of polynomial multiplication are the transformations  $a \rightarrow \hat{a}$  and the inverse  $\hat{a} \rightarrow a$  – these are the NTT and inverse NTT operations.

The other major time-consuming operation is the expansion of a seed  $\rho$  into the polynomial matrix  $\mathbf{A}$ . The matrix  $\mathbf{A}$  is needed for both signing and verification, therefore a good implementation of SHAKE-128 is important for the efficiency of the scheme.

The fastest NTT-implementation used in cryptography prior to this work is, to our knowledge, the AVX2 optimized NTT of NewHope [ADPS16] which uses floating point arithmetic. In our implementation, we resort to the more natural choice of using integer arithmetic, as for example also in [LN16]. Although we pack only 4 coefficients into one vector register of 256 bits, which is the same density that is also used by floating point implementations, we can improve on the multiplication speed by about a factor of 2 compared to the NewHope NTT adapted to the prime and degree used in Dilithium, see Table 1.2. We achieved this speed-up by carefully scheduling the instructions and interleaving the multiplications and reductions during the NTT so that parts of the multiplication latencies are hidden.<sup>1</sup> The effect of our fast NTT on the entire Dilithium scheme is a speed-up of 25% for signing and 15% for verification compared to the same implementation with the adapted floating point NTT from NewHope. We point out that in this implementation only the NTT and SHAKE are vectorized. By vectorizing other parts of the implementation like sampling and vector addition, a faster NTT would result in an even greater speed-up. Similarly, if one uses a faster algorithm for seed expansion of  $\rho$  into  $\mathbf{A}$  (or perhaps has  $\mathbf{A}$  already stored in memory if speed is truly of the essence), then the effect on the signing and (especially) verification algorithms will also be magnified.

<sup>1</sup>When the prime in the NTT is smaller than in Dilithium, additional improvements to the current algorithm produce a bigger speed-up (see Table 1.2 and [Sei18]).

**Table 1:** Haswell cycle counts on an Intel i7-4770 CPU in terms of  $10^3$  cycles of different AVX2 optimized NTT implementations for Dilithium and NewHope. Dilithium uses the ring  $\mathbb{Z}_q[X]/(X^{256} + 1)$  with  $q = 2^{23} - 2^{13} + 1$  and NewHope  $\mathbb{Z}_q[X]/(X^{1024} + 1)$  with  $q = 2^{13} + 2^{12} + 1$ . The cycle count for the forward NTT given in [ADPS16] is smaller than the cycle count of the inverse NTT because it is measured without the bitreversal operation. This operation can be left out if the input polynomial is randomly sampled with independent coefficients. In Dilithium this is not always the case, and so we compare the cycle counts including the bitreversal. The ADPS implementation for Dilithium is the floating point AVX2 implementation from NewHope [ADPS16] adapted to the Dilithium prime and degree. The LN implementation is the integer AVX2 implementation from [LN16]. Finally, the implementation for NewHope in the last row is similar to the one from this work but with major additional improvements that are not applicable to Dilithium, see [Sei18].

	Dilithium	NewHope
ADPS	3.2	9.5
LN	-	9.7
This work	1.5	2.2 [Sei18]

### 1.3 Related Work

The first asymptotically-efficient lattice-based signature scheme using the “Fiat-Shamir with Aborts” approach was constructed in [Lyu09] based on the Ring-SIS problem. The efficiency of the scheme was improved in [Lyu12] by basing the scheme on the combination of (Ring)-LWE and (Ring)-SIS problems. These schemes were further improved upon in [GLP12, BG14] by reducing the size of the signature. All the previous works had security proofs in the classical random oracle model.

Schemes based on the hardness of (Ring)-LWE in the quantum random oracle model were instantiated from “lossy identification” schemes [AFLT12] in [ABB<sup>+</sup>17] and [KLS17]. The latter work also showed that one could base the quantum security of all the schemes listed in the previous paragraph (as well as the scheme in this paper) on a non-interactive problem that’s a convolution of a lattice problem with a cryptographic hash function. The main disadvantage of schemes based only on (Ring)-LWE is that they are less efficient – they have substantially larger public key / signature sizes and also need to be instantiated over rings in which one cannot do NTT.

## 2 Basic Operations

### 2.1 Ring Operations

We let  $R$  and  $R_q$  respectively denote the rings  $\mathbb{Z}[X]/(X^n + 1)$  and  $\mathbb{Z}_q[X]/(X^n + 1)$ , for  $q$  an integer. Throughout this paper, the value of  $n$  will always be 256 and  $q$  will be the prime  $8380417 = 2^{23} - 2^{13} + 1$ . Regular font letters denote elements in  $R$  or  $R_q$  (which includes elements in  $\mathbb{Z}$  and  $\mathbb{Z}_q$ ) and bold lower-case letters represent column vectors with coefficients in  $R$  or  $R_q$ . By default, all vectors will be column vectors. Bold upper-case letters are matrices. For a vector  $\mathbf{v}$ , we denote by  $\mathbf{v}^T$  its transpose. The boolean operator  $\llbracket \text{statement} \rrbracket$  evaluates to 1 if statement is true, and to 0 otherwise.

**Modular reductions.** For an even (resp. odd) positive integer  $\alpha$ , we define  $r' = r \bmod^\pm \alpha$  to be the unique element  $r'$  in the range  $-\frac{\alpha}{2} < r' \leq \frac{\alpha}{2}$  (resp.  $-\frac{\alpha-1}{2} \leq r' \leq \frac{\alpha-1}{2}$ ) such that  $r' \equiv r \pmod{\alpha}$ . We will sometimes refer to this as a *centered* reduction modulo  $q$ . For any positive integer  $\alpha$ , we define  $r' = r \bmod^+ \alpha$  to be the unique element  $r'$  in the range  $0 \leq r' < \alpha$  such that  $r' \equiv r \pmod{\alpha}$ . When the exact representation is not important, we simply write  $r \bmod \alpha$ .

**Sizes of elements.** For an element  $w \in \mathbb{Z}_q$ , we write  $\|w\|_\infty$  to mean  $|w \bmod^\pm q|$ . We define the  $\ell_\infty$  and  $\ell_2$  norms for  $w = w_0 + w_1X + \dots + w_{n-1}X^{n-1} \in R$ :

$$\|w\|_\infty = \max_i \|w_i\|_\infty, \quad \|w\| = \sqrt{\|w_0\|_\infty^2 + \dots + \|w_{n-1}\|_\infty^2}.$$

Similarly, for  $\mathbf{w} = (w_1, \dots, w_k) \in R^k$ , we define

$$\|\mathbf{w}\|_\infty = \max_i \|w_i\|_\infty, \quad \|\mathbf{w}\| = \sqrt{\|w_1\|^2 + \dots + \|w_k\|^2}.$$

We will write  $S_\eta$  to denote all elements  $w \in R$  such that  $\|w\|_\infty \leq \eta$ .

## 2.2 NTT domain representation

Our modulus  $q$  is chosen such that there exists a 512-th root of unity  $r$  modulo  $q$ . Concretely, we always work with  $r = 1753$ . This implies that the cyclotomic polynomial  $X^{256} + 1$  splits into linear factors  $X - r^i$  modulo  $q$  with  $i = 1, 3, 5, \dots, 511$ . By the Chinese remainder theorem our cyclotomic ring  $R_q$  is thus isomorphic to the product of the rings  $\mathbb{Z}_q[X]/(X - r^i) \cong \mathbb{Z}_q$ . In this product of rings it is easy to multiply elements since the multiplication is pointwise there. The isomorphism

$$a \mapsto (a(r), a(r^3), \dots, a(r^{511})) : R_q \rightarrow \prod_i \mathbb{Z}_q[X]/(X - r^i)$$

can be computed quickly with the help of the Fast Fourier Transform. Since  $X^{256} + 1 = X^{256} - r^{256} = (X^{128} - r^{128})(X^{128} + r^{128})$  one can first compute the map

$$\mathbb{Z}_q[X]/(X^{256} + 1) \rightarrow \mathbb{Z}_q[X]/(X^{128} - r^{128}) \times \mathbb{Z}_q[X]/(X^{128} + r^{128})$$

and then continue separately with the two reduced polynomials of degree less than 128 noting that  $X^{128} + r^{128} = X^{128} - r^{384}$ . We give further detail about our NTT implementations in Section 5.1.

## 2.3 Hashing

Our scheme uses several different algorithms that hash strings in  $\{0, 1\}^*$  onto domains of various forms. Below we give the high level descriptions of these functions and defer the details of how exactly they are used in our signature scheme to Section 5.2.

**Hashing to a Ball.** Let  $B_h$  denote the set of elements of  $R$  that have  $h$  coefficients that are either  $-1$  or  $1$  and the rest are  $0$ . We have  $|B_h| = 2^h \cdot \binom{n}{h}$ . For our signature scheme, we will need a cryptographic hash function that hashes onto  $B_{60}$  (which has more than  $2^{256}$  elements). The algorithm we will use to create a random element in  $B_{60}$  is sometimes referred to as an “inside-out” version of the Fisher-Yates shuffle, and its high-level description is in Fig. 2.<sup>2</sup>

<sup>2</sup>Normally, the algorithm should begin at  $i = 0$ , but since there are 196 0's, the first 195 iterations would just be setting components of  $\mathbf{c}$  to 0.

```

SampleInBall
01 Initialize  $\mathbf{c} = c_0 c_1 \dots c_{255} = 00 \dots 0$ 
02 for  $i := 196$  to  $255$ 
03    $j \leftarrow \{0, 1, \dots, i\}$ 
04    $s \leftarrow \{0, 1\}$ 
05    $c_i := c_j$ 
06    $c_j := (-1)^s$ 
07 return  $\mathbf{c}$ 

```

**Figure 2:** Create a random 256-element array with 60  $\pm 1$ 's and 196 0's

**Expanding the Matrix  $\mathbf{A}$ .** The function `ExpandA` maps a uniform seed  $\rho \in \{0, 1\}^{256}$  to a matrix  $\mathbf{A} \in R_q^{k \times l}$  in CRT representation.

**Sampling the vectors  $y$ .** The function `ExpandMask`, used for deterministically generating the randomness of the signature scheme, maps  $K \parallel \mu \parallel \kappa$  to  $y \in S_{\gamma_1-1}^l$ .

**Collision resistant hash.** The function `CRH` used in our signature scheme is a collision resistant hash function mapping to  $\{0, 1\}^{384}$ .

## 2.4 High/Low Order Bits and Hints

To reduce the size of the public key, we will need some simple algorithms that extract “higher-order” and “lower-order” bits of elements in  $\mathbb{Z}_q$ . The goal is that when given an arbitrary element  $r \in \mathbb{Z}_q$  and another small element  $z \in \mathbb{Z}_q$ , we would like to be able to recover the higher order bits of  $r + z$  without needing to store  $z$ . We therefore define algorithms that take  $r, z$  and produce a 1-bit hint  $h$  that allows one to compute the higher order bits of  $r + z$  just using  $r$  and  $h$ . This hint is essentially the “carry” caused by  $z$  in the addition.

There are two different ways in which we will break up elements in  $\mathbb{Z}_q$  into their “high-order” bits and “low-order” bits. The first algorithm, `Power2Round $_q$` , is the straightforward bit-wise way to break up an element  $r = r_1 \cdot 2^d + r_0$  where  $r_0 = r \bmod^{\pm} 2^d$  and  $r_1 = (r - r_0)/2^d$ .

Notice that if we choose the representatives of  $r_1$  to be non-negative integers between 0 and  $\lfloor q/2^d \rfloor$ , then the distance (modulo  $q$ ) between any two  $r_1 \cdot 2^d$  and  $r'_1 \cdot 2^d$  is usually  $\geq 2^d$ , except for the border case. In particular, the distance modulo  $q$  between  $\lfloor q/2^d \rfloor \cdot 2^d$  and 0 could be very small. This is problematic in the case that we would like to produce a 1-bit hint, as adding a small number to  $r$  can actually cause the high-order bits of  $r$  to change by more than 1.

We avoid having the high-order bits change by more than 1 with a simple tweak. We select an  $\alpha$  that is a divisor of  $q - 1$  and write  $r = r_1 \cdot \alpha + r_0$  in the same way as before. For the sake of simplicity, we assume that  $\alpha$  is even (which is possible, as  $q$  is odd). The possible  $r_1 \cdot \alpha$ 's are now  $\{0, \alpha, 2\alpha, \dots, q - 1\}$ . Note that the distance between  $q - 1$  and 0 is 1, and so we remove  $q - 1$  from the set of possible  $r_1 \cdot \alpha$ 's, and simply round the corresponding  $r$ 's to 0. Because  $q - 1$  and 0 differ by 1, all this does is possibly increase the magnitude of the remainder  $r_0$  by 1. This procedure is called `Decompose $_q$` . Using this procedure as a sub-routine, we can define the `MakeHint $_q$`  and `UseHint $_q$`  routines that produce a hint and, respectively, use the hint to recover the high-order bits of the sum. For notational convenience, we also define `HighBits $_q$`  and `LowBits $_q$`  routines that simply extract  $r_1$  and  $r_0$ , respectively, from the output of `Decompose $_q$` .

The below Lemmas state the crucial properties of these supporting algorithms that are necessary for the correctness and security of our scheme. Their proofs can be found in



<p><u>Power2Round<sub>q</sub>(r, d)</u>  08 <math>r := r \bmod^+ q</math>  09 <math>r_0 := r \bmod^\pm 2^d</math>  10 <b>return</b> <math>((r - r_0)/2^d, r_0)</math></p> <p><u>MakeHint<sub>q</sub>(z, r, α)</u>  11 <math>r_1 := \text{HighBits}_q(r, \alpha)</math>  12 <math>v_1 := \text{HighBits}_q(r + z, \alpha)</math>  13 <b>return</b> <math>\llbracket r_1 \neq v_1 \rrbracket</math></p> <p><u>UseHint<sub>q</sub>(h, r, α)</u>  14 <math>m := (q - 1)/\alpha</math>  15 <math>(r_1, r_0) := \text{Decompose}_q(r, \alpha)</math>  16 <b>if</b> <math>h = 1</math> and <math>r_0 &gt; 0</math> <b>return</b> <math>(r_1 + 1) \bmod^+ m</math>  17 <b>if</b> <math>h = 1</math> and <math>r_0 \leq 0</math> <b>return</b> <math>(r_1 - 1) \bmod^+ m</math>  18 <b>return</b> <math>r_1</math></p>	<p><u>Decompose<sub>q</sub>(r, α)</u>  19 <math>r := r \bmod^+ q</math>  20 <math>r_0 := r \bmod^\pm \alpha</math>  21 <b>if</b> <math>r - r_0 = q - 1</math>  22     <b>then</b> <math>r_1 := 0; r_0 := r_0 - 1</math>  23 <b>else</b> <math>r_1 := (r - r_0)/\alpha</math>  24 <b>return</b> <math>(r_1, r_0)</math></p> <p><u>HighBits<sub>q</sub>(r, α)</u>  25 <math>(r_1, r_0) := \text{Decompose}_q(r, \alpha)</math>  26 <b>return</b> <math>r_1</math></p> <p><u>LowBits<sub>q</sub>(r, α)</u>  27 <math>(r_1, r_0) := \text{Decompose}_q(r, \alpha)</math>  28 <b>return</b> <math>r_0</math></p>
---	---

**Figure 3:** Supporting algorithms for Dilithium.

## Appendix A.

**Lemma 1.** *Suppose that  $q$  and  $\alpha$  are positive integers satisfying  $q > 2\alpha$ ,  $q \equiv 1 \pmod{\alpha}$  and  $\alpha$  even. Let  $\mathbf{r}$  and  $\mathbf{z}$  be vectors of elements in  $R_q$  where  $\|\mathbf{z}\|_\infty \leq \alpha/2$ , and let  $\mathbf{h}, \mathbf{h}'$  be vectors of bits. Then the  $\text{HighBits}_q$ ,  $\text{MakeHint}_q$ , and  $\text{UseHint}_q$  algorithms satisfy the following properties:*

1.  $\text{UseHint}_q(\text{MakeHint}_q(\mathbf{z}, \mathbf{r}, \alpha), \mathbf{r}, \alpha) = \text{HighBits}_q(\mathbf{r} + \mathbf{z}, \alpha)$ .
2. Let  $\mathbf{v}_1 = \text{UseHint}_q(\mathbf{h}, \mathbf{r}, \alpha)$ . Then  $\|\mathbf{r} - \mathbf{v}_1 \cdot \alpha\|_\infty \leq \alpha + 1$ . Furthermore, if the number of 1's in  $\mathbf{h}$  is  $\omega$ , then all except at most  $\omega$  coefficients of  $\mathbf{r} - \mathbf{v}_1 \cdot \alpha$  will have magnitude at most  $\alpha/2$  after centered reduction modulo  $q$ .
3. For any  $\mathbf{h}, \mathbf{h}'$ , if  $\text{UseHint}_q(\mathbf{h}, \mathbf{r}, \alpha) = \text{UseHint}_q(\mathbf{h}', \mathbf{r}, \alpha)$ , then  $\mathbf{h} = \mathbf{h}'$ .

**Lemma 2.** *If  $\|\mathbf{s}\|_\infty \leq \beta$  and  $\|\text{LowBits}_q(\mathbf{r}, \alpha)\|_\infty < \alpha/2 - \beta$ , then*

$$\text{HighBits}_q(\mathbf{r}, \alpha) = \text{HighBits}_q(\mathbf{r} + \mathbf{s}, \alpha).$$

## 3 Signature

The Key Generation, Signing, and Verification algorithms for our signature scheme are presented in Fig. 4. We present the deterministic version of the scheme in which the randomness used in the signing procedure is generated (using SHAKE-256) as a deterministic function of the message and a small secret key. Since our signing procedure may need to be repeated several times until a signature is produced, we also append a counter in order to make the SHAKE-256 output differ with each signing attempt of the same message. Also due to the fact that each message may require several iterations to sign, we compute an initial digest of the message using a collision-resistant hash function, and use this digest in place of the message throughout the signing procedure.

As discussed in Section 1.2, the main design improvement of Dilithium over the scheme in Fig. 1 is that the public key size is reduced by a factor of around 2.5 at the expense of an additional hundred bytes in the signature. To accomplish the size reduction, the key generation algorithm outputs  $\mathbf{t}_1 := \text{Power2Round}_q(\mathbf{t}, d)$  as the public key instead of  $\mathbf{t}$  as



```

Gen
01  $\rho \leftarrow \{0, 1\}^{256}$ 
02  $K \leftarrow \{0, 1\}^{256}$ 
03  $(\mathbf{s}_1, \mathbf{s}_2) \leftarrow S_\eta^\ell \times S_\eta^k$ 
04  $\mathbf{A} \in R_q^{k \times \ell} := \text{ExpandA}(\rho)$  //  $\mathbf{A}$  is stored in NTT Domain Representation
05  $\mathbf{t} := \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2$ 
06  $(\mathbf{t}_1, \mathbf{t}_0) := \text{Power2Round}_q(\mathbf{t}, d)$ 
07  $tr \in \{0, 1\}^{384} := \text{CRH}(\rho \parallel \mathbf{t}_1)$ 
08 return  $(pk = (\rho, \mathbf{t}_1), sk = (\rho, K, tr, \mathbf{s}_1, \mathbf{s}_2, \mathbf{t}_0))$ 

Sign $(sk, M)$ 
09  $\mathbf{A} \in R_q^{k \times \ell} := \text{ExpandA}(\rho)$  //  $\mathbf{A}$  is stored in NTT Domain Representation
10  $\mu \in \{0, 1\}^{384} := \text{CRH}(tr \parallel M)$ 
11  $\kappa := 0, (\mathbf{z}, \mathbf{h}) := \perp$ 
12 while  $(\mathbf{z}, \mathbf{h}) = \perp$  do
13    $\mathbf{y} \in S_{\gamma_1 - 1}^\ell := \text{ExpandMask}(K \parallel \mu \parallel \kappa)$ 
14    $\mathbf{w} := \mathbf{A}\mathbf{y}$ 
15    $\mathbf{w}_1 := \text{HighBits}_q(\mathbf{w}, 2\gamma_2)$ 
16    $c \in B_{60} := \text{H}(\mu \parallel \mathbf{w}_1)$ 
17    $\mathbf{z} := \mathbf{y} + c\mathbf{s}_1$ 
18    $(\mathbf{r}_1, \mathbf{r}_0) := \text{Decompose}_q(\mathbf{w} - c\mathbf{s}_2, 2\gamma_2)$ 
19   if  $\|\mathbf{z}\|_\infty \geq \gamma_1 - \beta$  or  $\|\mathbf{r}_0\|_\infty \geq \gamma_2 - \beta$  or  $\mathbf{r}_1 \neq \mathbf{w}_1$ , then  $(\mathbf{z}, \mathbf{h}) := \perp$ 
20   else
21      $\mathbf{h} := \text{MakeHint}_q(-c\mathbf{t}_0, \mathbf{w} - c\mathbf{s}_2 + c\mathbf{t}_0, 2\gamma_2)$ 
22     if  $\|c\mathbf{t}_0\|_\infty \geq \gamma_2$  or the # of 1's in  $\mathbf{h}$  is greater than  $\omega$ , then  $(\mathbf{z}, \mathbf{h}) := \perp$ 
23      $\kappa := \kappa + 1$ 
24 return  $\sigma = (\mathbf{z}, \mathbf{h}, c)$ 

Verify $(pk, M, \sigma = (\mathbf{z}, \mathbf{h}, c))$ 
25  $\mathbf{A} \in R_q^{k \times \ell} := \text{ExpandA}(\rho)$  //  $\mathbf{A}$  is stored in NTT Domain Representation
26  $\mu \in \{0, 1\}^{384} := \text{CRH}(\text{CRH}(\rho \parallel \mathbf{t}_1) \parallel M)$ 
27  $\mathbf{w}'_1 := \text{UseHint}_q(\mathbf{h}, \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d, 2\gamma_2)$ 
28 return  $\llbracket \|\mathbf{z}\|_\infty < \gamma_1 - \beta \rrbracket$  and  $\llbracket c = \text{H}(\mu \parallel \mathbf{w}'_1) \rrbracket$  and  $\llbracket \# \text{ of 1's in } \mathbf{h} \text{ is } \leq \omega \rrbracket$ 
    
```

**Figure 4:** Our signature scheme Dilithium.

in Fig. 1. This means that instead of  $\lceil \log q \rceil$  bits per coefficient, the public key requires  $\lceil \log q \rceil - d$  bits. In our instantiation,  $q \approx 2^{23}$  and  $d = 14$ , which means that instead of 23 bits in each public key coefficient, there are instead 9.

The main problem with not having the entire  $\mathbf{t}$  in the public key is that the verification algorithm is no longer able to exactly compute  $\mathbf{w}'_1$  in Line 13 in Fig. 1. In order to do this, the verification algorithm will need the high order bits of  $\mathbf{A}\mathbf{z} - c\mathbf{t}$ , but it can only compute  $\mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d = \mathbf{A}\mathbf{z} - c\mathbf{t} + c\mathbf{t}_0$ . But since the product  $c\mathbf{t}_0$  consists of only small numbers, and we only care about the high order bits, we really only need to know the carries that each coefficient  $c\mathbf{t}_0$  causes. These are the carries that the signer sends as a hint to the verifier. Heuristically, based on our parameter choices, there should not be more than  $\omega$  positions in which a carry is caused. The signer therefore simply sends the positions in which these carries occur (this is the extra bytes in the signature), which allows the verifier to compute the high order bits of  $\mathbf{A}\mathbf{z} - c\mathbf{t}$ .

### 3.1 Correctness

In this section, we prove the correctness of the signature scheme.

**Table 2:** Parameters for Dilithium.

	weak	medium	recommended	very high
$q$	8380417	8380417	8380417	8380417
$d$	14	14	14	14
weight of $c$	60	60	60	60
$\gamma_1 = (q - 1)/16$	523776	523776	523776	523776
$\gamma_2 = \gamma_1/2$	261888	261888	261888	261888
$(k, \ell)$	(3, 2)	(4, 3)	(5, 4)	(6, 5)
$\eta$	7	6	5	3
$\beta$	375	325	275	175
$\omega$	64	80	96	120
pk size (bytes)	896	1184	1472	1760
sig size (bytes)	1387	2044	2701	3366
Exp. reps (from Eq. (5))	4.3	5.9	6.6	4.3
BKZ block-size $b$ to break SIS	235	355	475	605
Best Known Classical bit-cost	68	103	138	176
Best Known Quantum bit-cost	62	94	125	160
Best Plausible bit-cost	48	73	98	125
BKZ block-size $b$ to break LWE	200	340	485	595
Best Known Classical bit-cost	58	100	141	174
Best Known Quantum bit-cost	53	91	128	158
Best Plausible bit-cost	41	71	100	124
Gen median cycles (Skylake)	149, 033	250, 774	377, 765	493, 926
Sign median cycles (Skylake)	658, 585	1, 157, 783	1, 769, 345	1, 566, 830
Sign average cycles (Skylake)	822, 232	1, 510, 580	2, 239, 468	2, 029, 811
Verify median cycles (Skylake)	174, 116	273, 570	390, 028	523, 443
Gen median cycles (AVX2, Skylake)	79, 421	128, 804	199, 486	258, 486
Sign median cycles (AVX2, Skylake)	213, 166	344, 174	507, 801	519, 776
Sign average cycles (AVX2, Skylake)	276, 851	448, 885	627, 120	636, 360
Verify median cycles (AVX2, Skylake)	81, 269	118, 704	174, 901	246, 116

If  $\|ct_0\|_\infty < \gamma_2$ , then by Lemma 1 we know that

$$\text{UseHint}_q(\mathbf{h}, \mathbf{w} - cs_2 + ct_0, 2\gamma_2) = \text{HighBits}_q(\mathbf{w} - cs_2, 2\gamma_2).$$

Since  $\mathbf{w} = \mathbf{A}\mathbf{y}$  and  $\mathbf{t} = \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2$ , we have that

$$\mathbf{w} - cs_2 = \mathbf{A}\mathbf{y} - cs_2 = \mathbf{A}(\mathbf{z} - cs_1) - cs_2 = \mathbf{A}\mathbf{z} - ct, \quad (2)$$

and  $\mathbf{w} - cs_2 + ct_0 = \mathbf{A}\mathbf{z} - ct_1 \cdot 2^d$ . Therefore the verifier computes

$$\text{UseHint}_q(\mathbf{h}, \mathbf{A}\mathbf{z} - ct_1 \cdot 2^d, 2\gamma_2) = \text{HighBits}_q(\mathbf{w} - cs_2, 2\gamma_2).$$

Furthermore, because the signer also checks in Line 19 that  $\mathbf{r}_1 = \mathbf{w}_1$ , this is equivalent to

$$\text{HighBits}_q(\mathbf{w} - cs_2, 2\gamma_2) = \text{HighBits}_q(\mathbf{w}, 2\gamma_2). \quad (3)$$

Therefore, the  $\mathbf{w}_1$  computed by the verifier is the same as that of the signer, and the verification procedure will always accept.

### 3.2 Number of Iterations

We now want to compute the probability that Step 19 will set  $(\mathbf{z}, \mathbf{h})$  to  $\perp$ . The probability that  $\|\mathbf{z}\|_\infty < \gamma_1 - \beta$  can be computed by considering each coefficient separately. For each coefficient  $\sigma$  of  $c\mathbf{s}_1$ , the corresponding coefficient of  $\mathbf{z}$  will be between  $-\gamma_1 + \beta + 1$  and  $\gamma_1 - \beta - 1$  (inclusively) whenever the corresponding coefficient of  $\mathbf{y}_i$  is between  $-\gamma_1 + \beta + 1 - \sigma$  and  $\gamma_1 - \beta - 1 - \sigma$ . The size of this range is  $2(\gamma_1 - \beta) - 1$ , and the coefficients of  $\mathbf{y}$  have  $2\gamma_1 - 1$  possibilities. Thus the probability that every coefficient of  $\mathbf{y}$  is in the good range is

$$\left(\frac{2(\gamma_1 - \beta) - 1}{2\gamma_1 - 1}\right)^{256 \cdot \ell} = \left(1 - \frac{\beta}{\gamma_1 - 1/2}\right)^{\ell n} \approx e^{-256 \cdot \beta \ell / \gamma_1}, \quad (4)$$

where we used the fact that our values of  $\gamma_1$  are large compared to  $1/2$ .

We now move to computing the probability that we have

$$\|\mathbf{r}_0\|_\infty = \|\text{LowBits}_q(\mathbf{w} - c\mathbf{s}_2, 2\gamma_2)\|_\infty < \gamma_2 - \beta.$$

If we (heuristically) assume that the low order bits are uniformly distributed modulo  $2\gamma_2$ , then there is a

$$\left(\frac{2(\gamma_2 - \beta) - 1}{2\gamma_2}\right)^{256 \cdot k} \approx e^{-256 \cdot \beta k / \gamma_2}$$

probability that all the coefficients are in the good range (using the fact that our values of  $\beta$  are large compared to  $1/2$ ).

As we already mentioned, if  $\|c\mathbf{s}_2\|_\infty \leq \beta$ , then  $\|\mathbf{r}_0\|_\infty < \gamma_2 - \beta$  implies that  $\mathbf{r}_1 = \mathbf{w}_1$ . Thus the last check should succeed with overwhelming probability when the previous check passed. Therefore, the probability that Step 19 passes is

$$\approx e^{-256 \cdot \beta (\ell / \gamma_1 + k / \gamma_2)}. \quad (5)$$

It is more difficult to formally compute the probability that Step 22 results in a restart. The parameters were set such that heuristically  $(\mathbf{z}, \mathbf{h}) = \perp$  with probability less than 1%. Therefore the vast majority of the loop repetitions will be caused by Step 19.

## 4 Security

The standard security notion for digital signatures is UF-CMA security, which is security under chosen message attacks. In this security model, the adversary gets the public key and has access to a signing oracle to sign messages of his choice. The adversary's goal is to come up with a valid signature of a new message. A slightly stronger security requirement that is sometimes useful is SUF-CMA (Strong Unforgeability under Chosen Message Attacks), which also allows the adversary to win by producing a different signature of a message that he has already seen.

It can be shown that in the (classical) random oracle model, Dilithium is SUF-CMA secure based on the hardness of the standard MLWE and MSIS lattice problems. The reduction, however, is not tight. Furthermore, since we also care about quantum attackers, we need to consider the security of the scheme when the adversary can query the hash function on a superposition of inputs (i.e. security in the quantum random oracle model – QROM). Since the classical security proof uses the “forking lemma” (which is essentially rewinding), the reduction does not transfer over to the quantum setting.

There are no counter-examples of schemes whose security is actually affected by the non-tightness of the reduction. For example, schemes like Schnorr signatures [Sch89], GQ signatures [GQ88], etc. all set their parameters ignoring the non-tightness of the reduction.

Furthermore, the only known uses of the additional power of quantum algorithms against schemes whose security is based on quantum-resistant problems under a classical reduction involve “Grover-type” algorithms that improve exhaustive search (although it has been shown that there cannot be a “black-box” proof that the Fiat-Shamir transform is secure in the QROM [ARU14]).

The reason that there haven’t been any attacks taking advantage of the non-tightness of the reduction is because there is an intermediate problem which is *tightly* equivalent, even under quantum reductions, to the UF-CMA security of the signature scheme. This problem is essentially a “convolution” of the underlying mathematical problem (such as MSIS or discrete log) with a cryptographic hash function  $H$ . It would appear that as long as there is no relationship between the structure of the math problem and  $H$ , solving this intermediate problem is not easier than solving the mathematical problem.<sup>3</sup>

Below, we will introduce the hardness assumptions upon whose hardness the SUF-CMA security of our scheme is based. The first two assumptions, MLWE and MSIS, are standard lattice problems which are a generalization of LWE, Ring-LWE, SIS, and Ring-SIS. The third problem, SelfTargetMSIS is the aforementioned problem that’s based on the combined hardness of MSIS and the hash function  $H$ . In the classical ROM, there is a (non-tight) reduction from MSIS to SelfTargetMSIS.

## 4.1 Assumptions

**The MLWE Problem.** For integers  $m, k$ , and a probability distribution  $D : R_q \rightarrow [0, 1]$ , we say that the advantage of algorithm  $A$  in solving the decisional  $\text{MLWE}_{m,k,D}$  problem over the ring  $R_q$  is

$$\begin{aligned} \text{Adv}_{m,k,D}^{\text{MLWE}} := & \left| \Pr[b = 1 \mid \mathbf{A} \leftarrow R_q^{m \times k}; \mathbf{t} \leftarrow R_q^m; b \leftarrow A(\mathbf{A}, \mathbf{t})] \right. \\ & \left. - \Pr[b = 1 \mid \mathbf{A} \leftarrow R_q^{m \times k}; \mathbf{s}_1 \leftarrow D^k; \mathbf{s}_2 \leftarrow D^m; b \leftarrow A(\mathbf{A}, \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2)] \right|. \end{aligned}$$

**The MSIS Problem.** To an algorithm  $A$  we associate the advantage function  $\text{Adv}_{m,k,\gamma}^{\text{MSIS}}$  to solve the (Hermite Normal Form)  $\text{MSIS}_{m,k,\gamma}$  problem over the ring  $R_q$  as

$$\text{Adv}_{m,k,\gamma}^{\text{MSIS}}(A) := \Pr \left[ 0 < \|\mathbf{y}\|_\infty \leq \gamma \wedge [\mathbf{I} \mid \mathbf{A}] \cdot \mathbf{y} = \mathbf{0} \mid \mathbf{A} \leftarrow R_q^{m \times k}; \mathbf{y} \leftarrow A(\mathbf{A}) \right].$$

**The SelfTargetMSIS Problem.** Suppose that  $H : \{0, 1\}^* \rightarrow B_{60}$  is a cryptographic hash function. To an algorithm  $A$  we associate the advantage function

$$\begin{aligned} \text{Adv}_{H,m,k,\gamma}^{\text{SelfTargetMSIS}}(A) := \\ \Pr \left[ \begin{array}{l} 0 \leq \|\mathbf{y}\|_\infty \leq \gamma \\ \wedge H([\mathbf{I} \mid \mathbf{A}] \cdot \mathbf{y} \parallel M) = c \end{array} \mid \mathbf{A} \leftarrow R_q^{m \times k}; \left( \mathbf{y} := \begin{bmatrix} \mathbf{r} \\ c \end{bmatrix}, M \right) \leftarrow A^{H(\cdot)}(\mathbf{A}) \right]. \end{aligned}$$

## 4.2 Signature Scheme Security

The concrete security of Dilithium was analyzed in [KLS17], where it was shown that if  $H$  is a quantum random oracle (i.e., a quantum-accessible perfect hash function), the advantage of an adversary  $A$  breaking the SUF-CMA security of the signature scheme is

$$\text{Adv}_{\text{Dilithium}}^{\text{SUF-CMA}}(A) \leq \text{Adv}_{k,\ell,D}^{\text{MLWE}}(B) + \text{Adv}_{H,k,\ell+1,\zeta}^{\text{SelfTargetMSIS}}(C) + \text{Adv}_{k,\ell,\zeta'}^{\text{MSIS}}(D) + 2^{-254},^4 \quad (6)$$

<sup>3</sup>In the ROM, there is indeed a (non-tight) reduction using the forking lemma that states that solving this problem is as hard as solving the underlying mathematical problem.

<sup>4</sup>To simplify the concrete security bound, we assume that  $\text{ExpandA}$  produces a uniform matrix  $\mathbf{A} \in R_q^{k \times \ell}$ ,  $\text{ExpandMask}(K, \cdot)$  is a perfect pseudo-random function, and  $\text{CRH}$  is a perfect collision-resistant hash function.

for  $D$  a uniform distribution over  $S_\eta$ , and

$$\zeta = \max\{\gamma_1 - \beta, 2\gamma_2 + 1 + 2^{d-1} \cdot 60\} \leq 4\gamma_2, \quad (7)$$

$$\zeta' = \max\{2(\gamma_1 - \beta), 4\gamma_2 + 2\} \leq 4\gamma_2 + 2. \quad (8)$$

Furthermore, if the running times and success probabilities (i.e. advantages) of  $A, B, C, D$  are  $t_A, t_B, t_C, t_D, \epsilon_A, \epsilon_B, \epsilon_C, \epsilon_D$ , then the lower bound on  $t_A/\epsilon_A$  is within a small multiplicative factor of  $\min t_i/\epsilon_i$  for  $i \in \{B, C, D\}$ .

Intuitively, the MLWE assumption is needed to protect against key-recovery, the SelfTargetMSIS is the assumption upon which new message forgery is based, and the MSIS assumption is needed for strong unforgeability. We will now sketch some parts of the security proof that are relevant to the concrete parameter setting.

#### 4.2.1 UF-CMA Security Sketch

It was shown in [KLS17] that for zero-knowledge deterministic signature schemes, if an adversary having quantum access to  $H$  and classical access to a signing oracle can produce a forgery of a new message, then there is also an adversary who can produce a forgery without access to the signing oracle (so he only gets the public key).<sup>5</sup> The latter security model is called UF-NMA – unforgeability under no-message attack. By the MLWE assumption, the public key  $(\mathbf{A}, \mathbf{t} = \mathbf{A}\mathbf{s}_1 + \mathbf{s}_2)$  is indistinguishable from  $(\mathbf{A}, \mathbf{t})$  where  $\mathbf{t}$  is chosen uniformly at random. The proof that our signature scheme is zero-knowledge is fairly standard and follows the framework from [Lyu09, Lyu12, BG14]. It is formally proved in [KLS17]<sup>6</sup> and we sketch the proof in Appendix B.

If we thus assume that the  $\text{MLWE}_{k,\ell,D}$  problem is hard, where  $D$  is the distribution that samples a uniform integer in the range  $[-\eta, \eta]$ , then to prove UF-NMA security, we only need to analyze the hardness of the experiment where the adversary receives a random  $(\mathbf{A}, \mathbf{t})$  and then needs to output a valid message/signature pair  $M, (\mathbf{z}, \mathbf{h}, c)$  such that

- $\|\mathbf{z}\|_\infty < \gamma_1 - \beta$
- $H(\text{UseHint}_q(\mathbf{h}, \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d, 2\gamma_2) \| M) = c$
- # of 1's in  $\mathbf{h}$  is  $\leq \omega$

Lemma 1 implies that one can rewrite

$$\text{UseHint}_q(\mathbf{h}, \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d, 2\gamma_2) = \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d + \mathbf{u}, \quad (9)$$

where  $\|\mathbf{u}\|_\infty \leq 2\gamma_2 + 1$ . Furthermore, only  $\omega$  coefficients of  $\mathbf{u}$  will have magnitude greater than  $\gamma_2$ . If we write  $\mathbf{t} = \mathbf{t}_1 \cdot 2^d + \mathbf{t}_0$  where  $\|\mathbf{t}_0\|_\infty \leq 2^{d-1}$ , then we can rewrite Eq. (9) as

$$\mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d + \mathbf{u} = \mathbf{A}\mathbf{z} - c(\mathbf{t} - \mathbf{t}_0) + \mathbf{u} = \mathbf{A}\mathbf{z} - c\mathbf{t} + (c\mathbf{t}_0 + \mathbf{u}) = \mathbf{A}\mathbf{z} - c\mathbf{t} + \mathbf{u}'. \quad (10)$$

Note that the worst-case upper-bound for  $\mathbf{u}'$  is

$$\|\mathbf{u}'\|_\infty \leq \|c\mathbf{t}_0\|_\infty + \|\mathbf{u}\|_\infty \leq \|c\|_1 \cdot \|\mathbf{t}_0\|_\infty + \|\mathbf{u}\|_\infty \leq 60 \cdot 2^{d-1} + 2\gamma_2 + 1 < 4\gamma_2.$$

<sup>5</sup>It was also shown in [KLS17] that the “deterministic” part of the requirement can be relaxed. The security proof simply loses a factor of the number of different signatures produced per message in its tightness. Thus, for example, if one were to implement the signature scheme (with the same secret key) on several devices with different random-number generators, the security of the scheme would not be affected much.

<sup>6</sup>In that paper, it is actually proved that the underlying zero-knowledge proof is zero-knowledge and then the security of the signature scheme follows via black box transformations.

Thus a (quantum) adversary who is successful at creating a forgery of a new message is able to find  $\mathbf{z}, c, \mathbf{u}', M$  such that  $\|\mathbf{z}\|_\infty < \gamma_1 - \beta$ ,  $\|c\|_\infty = 1$ ,  $\|\mathbf{u}'\|_\infty < 4\gamma_2$ , and  $M \in \{0, 1\}^*$  such that

$$\mathsf{H} \left( [\mathbf{A} \mid \mathbf{t} \mid \mathbf{I}_k] \cdot \begin{bmatrix} \mathbf{z} \\ c \\ \mathbf{u}' \end{bmatrix} \parallel M \right) = c. \quad (11)$$

Since  $(\mathbf{A}, \mathbf{t})$  is completely random, this is exactly the definition of the `SelfTargetMSIS` problem from above. A standard forking lemma argument can be used to show that an adversary solving the above problem in the (standard) random oracle model can be used to solve the `MSIS` problem. While giving a reduction using the forking lemma is a good “sanity check”, it is not particularly useful for setting parameters due to its lack of tightness. So how does one set parameters? The Fiat-Shamir transform has been used for over 3 decades (and we have been aware of the non-tightness of the forking lemma for two of them), yet the parameter settings for schemes employing it have ignored this loss in tightness. Implicitly, therefore, these schemes rely on the exact hardness of analogues (based on various assumptions such as discrete log [Sch89], one-wayness of RSA [GQ88], etc.) of the problem in Eq. (11).

The intuition for the security of the problem in Eq. (11) (and its discrete log, RSA, etc. analogues) is as follows: since  $\mathsf{H}$  is a cryptographic hash function whose structure is completely independent of the algebraic structure of its inputs, choosing some  $M$  “strategically” should not help – so the problem would be equally hard if the  $M$  were fixed. Then, again relying on the independence of  $\mathsf{H}$  and the algebraic structure of its inputs, the only approach for obtaining a solution appears to be picking some  $\mathbf{w}$ , computing  $\mathsf{H}(\mathbf{w} \parallel M) = c$ , and then finding  $\mathbf{z}, \mathbf{u}'$  such that  $\mathbf{A}\mathbf{z} + \mathbf{u}' = \mathbf{w} + c\mathbf{t}$ .<sup>7</sup> The hardness of finding such  $\mathbf{z}, \mathbf{u}'$  with  $\ell_\infty$ -norms less than  $4\gamma_2$  such that

$$\mathbf{A}\mathbf{z} + \mathbf{u}' = \mathbf{t}' \quad (12)$$

for some  $\mathbf{t}'$  is the problem whose concrete security we will be analyzing. Note that this is conservative because in Eq. (11)  $\|\mathbf{z}\|_\infty < \gamma_1 - \beta \approx 2\gamma_2$ . Furthermore, only  $\omega$  coefficients of  $\mathbf{u}'$  can be larger than  $2\gamma_2$ .

#### 4.2.2 The Addition of the Strong Unforgeability Property

To handle the strong-unforgeability requirement, one needs to handle an additional case. Intuitively, the reduction from `UF-CMA` to `UF-NMA` used the fact that a forgery of a new message will necessarily require the use of a challenge  $c$  for which the adversary has never seen a valid signature (i.e.,  $(\mathbf{z}, \mathbf{h}, c)$  was never an output by the signing oracle). To prove strong-unforgeability, we also have to consider the case where the adversary sees a signature  $(\mathbf{z}, \mathbf{h}, c)$  for  $M$  and then only changes  $(\mathbf{z}, \mathbf{h})$ . In other words, the adversary ends up with two valid signatures such that

$$\text{UseHint}_q(\mathbf{h}, \mathbf{A}\mathbf{z} - c\mathbf{t}_1 \cdot 2^d, 2\gamma_2) = \text{UseHint}_q(\mathbf{h}', \mathbf{A}\mathbf{z}' - c\mathbf{t}_1 \cdot 2^d, 2\gamma_2).$$

By Lemma 1, the above equality can be shown to imply that there exist  $\|\mathbf{z}\|_\infty \leq 2(\gamma_1 - \beta)$  and  $\|\mathbf{u}\|_\infty \leq 4\gamma_2 + 2$  such that  $\mathbf{A}\mathbf{z} + \mathbf{u} = \mathbf{0}$ .

### 4.3 Concrete Security Analysis

In Appendix C, we describe the best known lattice attacks against the problems in Eq. (6) upon which the security of our signature scheme is based. The best attacks involve finding short vectors in some lattice. The main difference between the `MLWE` and `MSIS`

<sup>7</sup>This is indeed the (non-tight) proof sketch in the classical random oracle model.

problems is that the MLWE problem involves finding a short vector in a lattice in which an “unusually short” vector exists. The MSIS problem, on the other hand, involves just finding a short vector in a random lattice. In knapsack terminology, the MLWE problem is a low-density knapsack, while MSIS is a high-density knapsack instance. The analysis for the two instances is slightly different and we analyze the MLWE problem in Appendix C.2 and the MSIS problem (as well as SelfTargetMSIS) in Appendix C.3.

We follow the general methodology from [ADPS16, BCD<sup>+</sup>16] to analyze the security of our signature scheme, with minor adaptations. This methodology is significantly more conservative than prior ones used in lattice-based cryptography. In particular, we assume the adversary can run the asymptotically best algorithms known, with no overhead compared to the asymptotic run-times. In particular, we assume the adversary can cheaply handle huge amounts of (possibly quantum) memory. This conservatism is in line with the goal of long-term post-quantum security. We note that despite this security analysis methodology, our schemes remain competitive in practice.

While the MLWE and MSIS problems are defined over polynomial rings, we do not currently have any way of exploiting this ring structure, and therefore the best attacks are mounted by simply viewing these problems as LWE and SIS problems. The LWE and SIS problems are exactly as in the definitions of MLWE and MSIS in Section 4.1 with the ring  $R_q$  being replaced by  $\mathbb{Z}_q$ .

## 5 Implementation Details

### 5.1 NTT domain representation

Natural fast NTT implementations do not output vectors with coefficients in the order  $a(r), a(r^3), \dots, a(r^{511})$ . Therefore we define the NTT domain representation  $\hat{a} \in \mathbb{Z}_q^{256}$  of a polynomial  $a \in R_q$  to have coefficients in the order as output by our reference NTT. Concretely,

$$\hat{a} = (a(r_0), a(-r_0), \dots, a(r_{127}), a(-r_{127}))$$

where  $r_i = r^{\text{brv}(128+i)}$  with  $\text{brv}(k)$  the bitreversal of the 8 bit number  $k$ .

This is important for sampling the matrix  $\mathbf{A}$ . The matrix  $\mathbf{A}$  is only needed for multiplication. Hence, for the sake of faster implementations, the expansion function `ExpandA` does not output  $A \in R_q^{k \times l}$ . Instead it outputs a matrix with coefficients in  $\mathbb{Z}_q^{256}$ , which is interpreted as the NTT domain representation of  $\mathbf{A}$ . As  $\mathbf{A}$  needs to be sampled uniformly and the NTT is an isomorphism, `ExpandA` also needs to sample uniformly in this representation. To be compatible to Dilithium, an implementation whose NTT produces differently ordered vectors than our reference NTT needs to sample coefficients in a non-consecutive order.

### 5.2 Hashing

**Hashing to a Ball.** We now precisely describe the operation of the function  $H : \mu \parallel \mathbf{w}_1 \mapsto c \in B_{60}$  described in Fig. 2 as it is used in our signature scheme. The vector  $\mathbf{w}_1$  consists of  $k$  polynomials  $w_{1,0}, \dots, w_{1,k-1}$  in  $R_q$  with coefficients in  $\{0, \dots, 15\}$ . This allows  $\mathbf{w}_1$  to be packed in a string of  $k \cdot 256 \cdot 4/8 = k \cdot 128$  bytes. Concretely, the bytes numbers  $128i$  to  $128i + 127$ ,  $i = 0, \dots, k - 1$ , encode the coefficients of  $w_{1,i}$  in increasing order, where each byte encodes two consecutive coefficients of  $w_{1,i}$  in its low 4 bits and high 4 bits, respectively. For example, the first two coefficients  $c_0$  and  $c_1$  of  $w_{1,i}$  are encoded in the  $128i^{\text{th}}$  byte of the bit-packed string of  $\mathbf{w}_1$  as  $c_1 \cdot 16 + c_0$ .

H then absorbs the 384-bit string  $\mu$  immediately followed by the  $k \cdot 128$  bytes for  $\mathbf{w}_1$  into SHAKE-256. Throughout its operations the function squeezes SHAKE-256 in order to obtain a stream of random bytes of variable length. The first 60 bits in the first 8 bytes



of this random stream are interpreted as 60 random sign bits  $s_i \in \{0, 1\}$ . The remaining 4 bits are discarded. Then H uses Algorithm 2 to compute  $c$  where in each iteration of the for loop it uses rejection sampling on elements from  $\{0, \dots, 255\}$  until it gets a  $j \in \{0, \dots, i\}$ . An element in  $\{0, \dots, 255\}$  is obtained by interpreting the next byte of the random stream from SHAKE-256 as a number in this set. For the sign  $s$  the corresponding  $s_{i-196}$  is used.

**Expanding the Matrix  $\mathbf{A}$ .** The function `ExpandA` maps a uniform seed  $\rho \in \{0, 1\}^{256}$  to a matrix  $\mathbf{A} \in R_q^{k \times l}$  in NTT domain representation. It computes each coefficient  $\hat{a}_{i,j} \in \mathbb{Z}_q^{256}$  of  $\mathbf{A}$  separately. For the coefficient  $\hat{a}_{i,j}$  it absorbs the 32 bytes of  $\rho$  immediately followed by one byte representing  $0 \leq 16 \cdot j + i < 255$  into SHAKE-128. Next it uses consecutive blocks of 3 bytes of the variable-length output string in order to obtain a sequence of integers between 0 and  $2^{23} - 1$ . This is done by setting the highest bit of the third byte in each block to zero and interpreting the blocks in little endian byte order. So for example the three bytes  $b_0, b_1$  and  $b_2$  from SHAKE-128 are used to get the integer  $0 \leq b'_2 \cdot 2^{16} + b_1 \cdot 2^8 + b_0 \leq 2^{23} - 1$  where  $b'_2$  is the logical AND of  $b_2$  and  $2^{128} - 1$ . Finally, `ExpandA` performs rejection sampling on these 23-bit integers to sample the 256 coefficients  $a_{i,j}(r_0), a_{i,j}(-r_0), \dots, a_{i,j}(r_{127}), a_{i,j}(-r_{127})$  of  $\hat{a}_{i,j}$  uniformly from the set  $\{0, \dots, q - 1\}$  in the order of our NTT domain representation.

**Sampling the vectors  $y$ .** The function `ExpandMask` maps  $K \parallel \mu \parallel \kappa$  to  $y \in S_{\gamma_1-1}^l$ , where  $\kappa \geq 0$ , and works as follows. It computes each of the  $l$  coefficients of  $y$ , which are polynomials in  $S_{\gamma_1-1}$ , independently. For the  $i$ -th polynomial,  $0 \leq i < l$ , it absorbs 32 bytes of  $K$  concatenated with the 48 bytes of  $\mu$  and two bytes representing  $l\kappa + i$  in little endian byte order into SHAKE-256. Then each block of 5 consecutive output bytes is used to get two 20 bit integers between 0 and  $2^{20} - 1$ . For this the first two bytes of each output block together with a third byte having as lower 4 bits the lower 4 bits of the third output byte and 4 high zero bits is interpreted in little endian order. Then the high 4 bits of the third output byte followed by the 16 bits of the fourth and fifth byte are interpreted as the second 20 bit integer. As an example assume we have received the five bytes  $b_0, \dots, b_4$  from SHAKE-128. Then `ExpandMask` computes the two integers  $0 \leq b'_2 \cdot 2^{16} + b_1 \cdot 2^8 + b_0 \leq 2^{20} - 1$  and  $0 \leq b_4 \cdot 2^{12} + b_3 \cdot 2^4 + b'_2 \leq 2^{20} - 1$  where  $b'_2$  is the AND of  $b_2$  and 15 and  $b'_2 = \lfloor b_2/16 \rfloor$ . On the resulting sequence of 20 bit integers rejection sampling is performed to get 256 values  $v_j \in \{0, \dots, 2\gamma_1 - 2\}$ . From these the polynomial coefficients are computed in increasing order as  $q + \gamma_1 - 1 - v_j$ .

**Collision resistant hash.** The function CRH in Figure 4 is a collision resistant hash function. For this purpose 384 bits of the output of SHAKE-256 are used. CRH is called with two different sets of inputs. First it is called with  $\rho \parallel \mathbf{t}_1$ . The function then absorbs the 32 bytes of  $\rho$  followed by the  $k \cdot 256 \cdot 9/8$  bytes for the bit-packed representation of  $\mathbf{t}_1$  into SHAKE-256 and takes the first 48 bytes of the first output block of SHAKE-256 as the output hash. The bit-packing for  $\mathbf{t}_1$  is done in the same way as described above for  $\mathbf{w}_1$ . As always integers are read in little endian byte order. The second input is  $\mu \parallel M$ . Here the concatenation of the hash  $\mu$  and the message string are absorbed into SHAKE-256 and the first 48 output bytes are used as the resulting hash.

### 5.3 Data layout of keys and signature

**Public key.** The public key, containing  $\rho$  and  $\mathbf{t}_1$ , is stored as the concatenation of the bit-packed representations of  $\rho$  and  $\mathbf{t}_1$  in this order. Therefore, it has a size of  $32 + 288k$  bytes.

**Secret key.** The secret key contains  $\rho$ ,  $K$ ,  $tr$ ,  $\mathbf{s}_1$ ,  $\mathbf{s}_2$  and  $\mathbf{t}_0$  and is also stored as a bit-packed representation of these quantities in the given order. We explain the bit-packing of  $\mathbf{s}_1$ ,  $\mathbf{s}_2$  and  $\mathbf{t}_0$ . The polynomials inside  $\mathbf{s}_1$  and  $\mathbf{s}_2$  have coefficients with infinity norm at most  $\eta$ . So every coefficient  $c \in \mathbb{Z}_q$  of these polynomials can be represented by  $c = q + \eta - v$  with some  $v \in \{0, \dots, 2\eta\}$ . In the bit packing the values  $v$  are stored so that each polynomial needs  $256 \cdot 4/8$  bytes for the weak, medium and recommended security levels, and  $256 \cdot 3/8$  bytes for the very high security level. The coefficients of the polynomials of  $\mathbf{t}_0$  can be written in the form  $q + 2^{d-1} - v$  with  $v \in \{0, \dots, 2^d - 1\}$  and the representation by these  $v$  is used in the bit packing. Hence, this results in  $256 \cdot d/8$  bytes per polynomial. Consequently, a secret key requires  $64 + 48 + 32((k+l) \cdot \lceil \log(2\eta+1) \rceil + 14k)$  bytes. For the weak, medium and high security level this is equal to  $112 + 576k + 128l$  bytes. With the very high security parameters one needs  $112 + 544k + 96l = 3856$  bytes.

**Signature.** The signature bit string is the concatenation of a bit packed representation of  $\mathbf{z}$  and encodings of  $h$  and  $c$  in this order. The coefficients of the polynomials of  $\mathbf{z}$  are written in the form  $q + \gamma_1 - 1 - v$  with  $v \in \{0, \dots, 2\gamma_1 - 2\}$  and these values  $v$  are bit packed resulting in  $256 \cdot 20/8$  bytes per polynomial. Next we describe the encoding of  $h$ , which needs  $\omega + K$  bytes. Together all the polynomials in the vector  $\mathbf{h}$  have at most  $\omega$  non-zero coefficients. It is sufficient to store the locations of these non-zero coefficients. Each of the first  $\omega$  bytes of the bit string representing  $\mathbf{h}$  is the index  $i$  of the next non-zero coefficient in its polynomial, i.e.  $0 \leq i \leq 255$ , or zero if there are no more non-zero coefficients. The bytes numbers  $\omega$  up to  $\omega + k - 1$  record the  $k$  positions  $j$  of the polynomial boundaries in the string of  $\omega$  coefficient indices, where  $0 \leq j \leq \omega$ . In the encoding of the challenge  $c$ , the first 256 bits are 0 or 1 when the corresponding coefficient of  $c$  is zero or non-zero, respectively. The next 60 bits are 0 or 1 if the corresponding non-zero coefficient is 1 or  $-1$ , respectively. Note that there are precisely 60 non-zero coefficients. The 4 bits up to the next byte boundary are zero. Therefore, a signature requires  $640l + \omega + k + 40$  bytes.

## 5.4 Constant time implementation

Our reference implementation does not branch depending on secret data and does not access memory locations that depend on secret data. For the modular reductions that are needed for the arithmetic in  $R_q$  we never use the '%' operator of the C programming language. Instead we use Montgomery reductions without the correction steps and special reduction routines that are specific to our modulus  $q$ . For computing the rounding functions described in Section 2.4, we have implemented branching-free algorithms. On the other hand, when it is safe to reveal information, we have not tried to make the code constant-time. This includes the computation of the challenges and the rejection conditions in the signing algorithm. When performing rejection sampling, our code reveals which of the conditions was the reason for the rejection, and in case of the norm checks, which coefficient violated the bound. This is safe since the rejection probabilities for each coefficient are independent of secret data. The challenges reveal information about  $\text{CRH}(\mu \parallel \mathbf{w}_1)$  also in the case of rejected  $\mathbf{y}$ , but this does not reveal any information about the secret key when CRH is modeled as a random oracle and  $\mathbf{w}_1$  has high min-entropy.

## 5.5 Reference implementation

Our reference NTT is a natural iterative implementation for 32 bit unsigned integers that uses Cooley-Tukey butterflies in the forward transform and Gentleman-Sande butterflies in the inverse transform. For modular reductions after multiplying with a precomputed root of unity we use the Montgomery algorithm as was already done before in e.g. [ADPS16]. In order that the reduced values are correct representatives, the precomputed roots contain the Montgomery factor  $2^{32} \bmod q$ . We also use Montgomery reductions after the pointwise

product of the polynomials in the NTT domain representations. Since we cannot get the Montgomery factor in at this point, these products are in fact Hensel remainders  $r' \equiv r2^{32} \pmod{q}$ . We then make use of the fact that the NTT transform is linear and multiply by an additional Montgomery factor after the inverse NTT when we divide out the factor 256.

The implementations of the functions `ExpandA` and `ExpandMask` initially squeeze a number of output blocks of SHAKE-256 and SHAKE-128 that gives enough randomness with high probability. In the case of `ExpandA`, which samples uniform polynomials and hence needs at least  $3 \cdot 256 = 768$  random bytes per polynomial, 5 blocks from SHAKE-128 of 168 bytes each are needed at least for one polynomial. They suffice with probability greater than  $1 - 2^{-132}$ . `ExpandMask` initially retrieves 5 blocks from SHAKE-256 that have 136 bytes. This is the minimum number of blocks and suffices with probability greater than  $1 - 2^{-81}$ .

As mentioned in the introduction our reference implementation is protected against timing attacks. For this reason the centralized remainders in the rounding functions given in Figure 3 are not computed with branchings. Instead we use the following well-known trick to compute the centralized remainder  $r' = r \bmod^{\pm} \alpha$  where  $0 \leq r \leq 3\alpha/2$ . Subtracting  $\alpha/2 + 1$  from  $r$  yields a negative result if and only if  $r \leq \alpha/2$ . Therefore, shifting this result arithmetically to the right by 31 bits gives  $-1$ , i.e. the integer with all bits equal to 1, if  $r \leq \alpha/2$  and 0 otherwise. Then the logical AND of the shifted value and  $\alpha$  is added to  $r$  and  $\alpha/2 - 1$  subtracted. This results in  $r - \alpha$  if  $r > \alpha/2$  and  $r$  if  $r \leq \alpha/2$ , i.e. the centralized remainder.

We make heavy use of lazy reduction in our implementation. In the NTT we do not reduce the results of additions and subtractions at all. For rounding and norm checking it is important to map to standard representatives. This freezing of the coefficients is achieved in constant-time by conditionally subtracting  $q$  with another instance of the arithmetic right shift trick.

## 5.6 AVX2 optimized implementation

We have written an optimized implementation of Dilithium for CPUs that supports the AVX2 instruction set. Since the two most time-consuming operations are polynomial multiplication and the expansion of the matrix and vectors, the optimized implementation speeds up these two operations.

For polynomial multiplication, we have implemented a vectorized version of the NTT. This NTT achieves a full multiplication of two polynomials including three NTTs and the pointwise multiplication in less than 5000 Haswell cycles and is about a factor of 4.5 faster than the reference C code compiled using gcc with full machine-specific optimizations turned on. We do not use floating point instructions, where modular reductions are easily done by multiplying with a floating point inverse of  $q$  and rounding to get the quotient from which the remainder can be computed with another multiplication and a subtraction. Instead, we use integer instructions only and the same Montgomery reduction methodology as in the reference C code.

At any time our AVX2 optimized NTT has 32 unsigned integer coefficients, of 32 bits each, loaded into 8 AVX2 vector registers. Each of these vector registers then contains 4 extended 64 bit coefficients. So after three levels of NTT the reduced polynomials fit completely into these 8 registers and we can transform them to linear factors without further loads and stores. In the second last and last level the polynomials have degree less than 4. This means that every polynomial fits into one register but only half of the coefficients need to be multiplied by roots. For this reason we shuffle the vectors in order to group together coefficients that need to be multiplied. The instruction that we use for this task are `perm2i128` in the second last level and a combination of `vpsufd` and `vblendd` in the last level. The multiplications with the constant roots of unity are performed using the `vpmuldq` instruction. This instruction computes a full 64 bit product

of two 32 bit integers. It has a latency of 5 cycles on both Haswell and Skylake. In each level of the NTT half of the coefficients need to be multiplied. Therefore we can do four vector multiplications and Montgomery reductions in parallel. This hides some of the latency of the multiplication instructions.

For faster matrix and vector expansion, we use a vectorized SHAKE implementation that operates on 4 parallel sponges and hence can absorb and squeeze blocks in and out of these 4 sponges at the same time. For sampling this means that up to four coefficients can be sampled simultaneously.

## 5.7 Computational Efficiency

We have performed timing experiments with our reference implementation on a Skylake CPU. The results are presented in Table 2. They include the number of CPU cycles needed by the three operations key generation, signing and signature verification. These numbers are the medians of 10000 executions each. For the signing operation we also give the average number of cycles because the use of rejection sampling during signing has the effect that sometimes more time is needed than is indicated by the median cycle counts. Signing was performed with a message size of 32 bytes. The computer we have used is equipped with an Intel Core i7-6600U CPU. The system runs Ubuntu 18.4 with Linux Kernel version 4.15.0 and the code was compiled with gcc 7.3.0.

## 6 Comparisons

There are three major approaches for post-quantum signatures: hash-based signatures, multivariate signatures, and lattice-based signatures. Applications that can securely maintain a state (i.e., a constantly changing secret key) will certainly want to adopt stateful hash-based signatures like XMSS [BDH11], which is currently being standardized by the crypto forum research group (CFRG) of IETF. However, many applications are unlikely to migrate to a signature scheme that becomes insecure if, for example, the secret key is part of a backup and is restored to an older state. In their call for proposals, NIST therefore explicitly asks for stateless signature proposals. We present a comparison of our digital signature to other lattice-based and non-lattice-based schemes in Table 3.

**Non-Lattice-Based Signatures.** Arguably the most conservative approach to digital signature construction is using stateless hash-based signatures that come with tight reductions in the standard model to standard properties of a cryptographic hash function, like collision resistance or second-preimage resistance. Unfortunately, stateless hash-based signatures are rather inefficient in terms of signature size and signing speed. For example, the SPHINCS signature scheme has signatures of about 40 KB and takes about 50 million cycles to sign even on large Intel processors with the AVX2 vector instruction set.

Applications that need small signatures may want to turn their attention to multivariate signature schemes like Rainbow or HFEv-. These small signatures, however, come at the expense of rather large public keys. Many applications, for instance TLS, need to transmit public keys almost as much as signatures, so what becomes important is the sum of the public-key size and signature size. Also, efficient multivariate signature schemes do not enjoy a reduction from the  $\mathcal{MQ}$  problem. Instead these schemes require a hidden structure in the instance of the  $\mathcal{MQ}$  problem, which was a weakness exploited by many attacks on prior protocols (c.f. [KS98, DFSS07, FGP<sup>+</sup>15, TW12]). At Asiacrypt 2016, Chen, Hülsing, Rijneveld, Samardjiska, and Schwabe presented MQDSS [CHR<sup>+</sup>16], a multivariate signature scheme that *does* have a reduction from  $\mathcal{MQ}$ , but this reduction comes at the price of signatures that are not shorter than SPHINCS signatures.

**Lattice-Based Signatures.** The most efficient (in terms of key and signature sizes) lattice-based schemes are those based on the hardness of the NTRU problem. The one with the smallest key/signature combination [DLP14] uses the key-generation of [HHGP<sup>+</sup>03] to create lattices which are optimal for using the secure signing algorithm in [GPV08]. The “Fiat-Shamir with Aborts” based scheme BLISS [DDLL13] is also based on the hardness of finding short vectors in NTRU lattices and has slightly larger parameters.

While the above schemes have the smallest outputs of all lattice-based schemes, they contain several downsides. The most important one is that the schemes crucially require sampling from the discrete Gaussian distribution, and to the best of our knowledge no constant-time implementations of these schemes exist due to this “feature”.<sup>8</sup> A second downside is that the security of these schemes is based on NTRU rather than Ring/Module-LWE. The geometric structure of NTRU lattices has been recently exploited [KF17] to produce significantly better attacks against large-modulus/small-secret version of the problem (these attacks, though, currently do not extend to the parameters used in digital signatures). The final downside is that increasing/decreasing the security levels would require a complete re-implementation of the schemes. In particular, it is not possible to efficiently instantiate them in a way that constructs the public key out of small blocks, as in this paper.

At the other extreme of lattice constructions are digital signatures based on the hardness of standard lattice problems without any algebraic structure. The main downside of these schemes (c.f. [Lyu12, ABB<sup>+</sup>17]) is that they have extremely large public keys, and are not suitable for many practical applications.

Carefully optimized signatures based on the hardness of ideal or module lattice problems sit in the sweet spot that offers reasonably small signatures and public keys, good signing and verification speeds, and simple constant-time implementations. In particular, the combined size of the public key and signature of the scheme in the current paper is smaller than in all the non-lattice-based schemes that we are aware of. Starting from [Lyu09], there have been many improvements and implementations of this type of scheme (e.g. [Lyu12, GLP12, GOPS13, BG14]). The most practical one of these that existed prior to the current work is [BG14], which is essentially the scheme in Fig. 1. Since Dilithium reduces the public key of that scheme by a factor  $\approx 2.5$  (saving over 2KB) at the expense of an additional 100 bytes in the signature, it is the most efficient of all those contained in this line of work. We should mention that it is also possible to produce slightly-shorter parameters (around 10 - 15% saving in the signature size) if one samples from the Gaussian distribution, but we believe that this saving does not justify the added complexity of securely implementing the scheme.

---

<sup>8</sup>A recent work [MW17] addresses some issues, but their implementation still uses lookup tables which are not protected against side-channel attacks.

**Table 3:** Comparison of post-quantum signature schemes. Benchmarks were performed on an Intel Core i7-6600U (Skylake) if not indicated otherwise. The “sec” column states the security as reported by the authors of the respective papers. Cycles are stated for signing (**S**) and verification (**V**). Bytes are given for public keys (**pk**), and signatures (**s**). The column “ct?” indicates whether the software is protected against timing attacks.

Scheme	sec	ct?		Cycles		Bytes
<b>NTRU-based lattice signature schemes</b>						
NTRU-GPV [DLP14] (adapted to a 1024-dim instance) <sup>a</sup>	≫ 128	no	<b>S:</b>	??	<b>pk:</b>	1 792
			<b>V:</b>	??	<b>s:</b>	≈ 1 200
BLISS [DDLL13] (adapted to a 1024-dim instance) <sup>a</sup>	≫ 128	no	<b>S:</b>	??	<b>pk:</b>	1 792
			<b>V:</b>	??	<b>s:</b>	≈ 1 700
<b>Ring/Module-lattice-Based signature schemes</b>						
Dilithium, recomm. param. (this paper)	125	yes	<b>S:</b>	508K	<b>pk:</b>	1 472
			<b>V:</b>	175K	<b>s:</b>	2 701
<b>Standard-lattice-based signature schemes</b>						
TESLA-I [ABB <sup>+</sup> 17]	98	yes	<b>S:</b>	143 402K <sup>d</sup>	<b>pk:</b>	12MB
			<b>V:</b>	19 284K <sup>d</sup>	<b>s:</b>	2 444
<b>Stateless hash-based signature schemes</b>						
SPHINCS [BHH <sup>+</sup> 15]	128	yes	<b>S:</b>	51 636K <sup>b</sup>	<b>pk:</b>	1 056
			<b>V:</b>	1 451K <sup>b</sup>	<b>s:</b>	41 000
<b>Multivariate signature schemes</b>						
HmFEv(256,15,3,16) [CLP <sup>+</sup> 17]	128	yes	<b>S:</b>	1 497K <sup>c</sup>	<b>pk:</b>	83 100
			<b>V:</b>	15K <sup>c</sup>	<b>s:</b>	61
Rainbow(16,32,32,32) [CLP <sup>+</sup> 17]	128	yes	<b>S:</b>	68K <sup>c</sup>	<b>pk:</b>	145 500
			<b>V:</b>	22K <sup>c</sup>	<b>s:</b>	48
MQDSS [CHR <sup>+</sup> 16]	128	yes	<b>S:</b>	8 510K <sup>d</sup>	<b>pk:</b>	72
			<b>V:</b>	5 752K <sup>d</sup>	<b>s:</b>	40 952

<sup>a</sup> The scheme in the paper used a different security analysis and therefore considered instances of dimensions (i.e. 512) that are not secure enough using the security analysis in the current paper. Doubling the dimension to 1024 makes these schemes well over 128-bit quantum-secure. We adjusted the approximate parameters to match this increase in the instance size. We should point out that it might also be possible to increase the dimension to less than 1024 by not using a ring of the form  $\mathbb{Z}[X]/(X^N + 1)$ . This may reduce the parameters of the schemes.

<sup>b</sup> Benchmarked on an Intel Xeon E3-1275 (Haswell)

<sup>c</sup> Benchmarked on an Intel Xeon E3-1245 v3 (Haswell)

<sup>d</sup> Benchmarked on an Intel Core i7-4770K (Haswell)



## 7 Acknowledgements

We are very thankful to the referees for their useful suggestions. This work is supported by the European Commission through the ICT program under contracts ICT-645622 (PQCRYPTO), ICT-644729 (SAFEcrypto), and through the ERC Starting Grant ERC-2013-StG-335086 (LATTAC). It is also supported by the Swiss National Science Foundation through the 2014 transfer ERC Starting Grant CRETP2-166734 (FELICITY) and by the Netherlands Organization for Scientific Research (NWO) through Veni grant 639.021.645 (Cryptanalysis of Lattice-based Cryptography).

## References

- [ABB<sup>+</sup>17] Erdem Alkim, Nina Bindel, Johannes A. Buchmann, Özgür Dagdelen, Edward Eaton, Gus Gutoski, Juliane Krämer, and Filip Pawlega. Revisiting TESLA in the quantum random oracle model. In *PQCrypto*, pages 143–162, 2017. 5, 20, 21
- [ADPS16] Erdem Alkim, Léo Ducas, Thomas Pöppelmann, and Peter Schwabe. Post-quantum key exchange – a new hope. In *Proceedings of the 25th USENIX Security Symposium*, pages 327–343. USENIX Association, 2016. <http://cryptojedi.org/papers/#newhope>. 4, 5, 15, 17, 29
- [AFLT12] Michel Abdalla, Pierre-Alain Fouque, Vadim Lyubashevsky, and Mehdi Tibouchi. Tightly-secure signatures from lossy identification schemes. In *EUROCRYPT*, pages 572–590, 2012. 5
- [AG11] Sanjeev Arora and Rong Ge. New algorithms for learning in presence of errors. In Luca Aceto, Monika Henzinger, and Jiri Sgall, editors, *ICALP 2011, Part I*, volume 6755 of *LNCS*, pages 403–415. Springer, Heidelberg, July 2011. 30
- [ARU14] Andris Ambainis, Ansis Rosmanis, and Dominique Unruh. Quantum attacks on classical proof systems: The hardness of quantum rewinding. In *FOCS*, pages 474–483, 2014. 12
- [BCD<sup>+</sup>16] Joppe W. Bos, Craig Costello, Léo Ducas, Ilya Mironov, Michael Naehrig, Valeria Nikolaenko, Ananth Raghunathan, and Douglas Stebila. Frodo: Take off the ring! Practical, quantum-secure key exchange from LWE. In Edgar R. Weippl, Stefan Katzenbeisser, Christopher Kruegel, Andrew C. Myers, and Shai Halevi, editors, *ACM CCS 16*, pages 1006–1018. ACM Press, October 2016. 15
- [BDGL16] Anja Becker, Léo Ducas, Nicolas Gama, and Thijs Laarhoven. New directions in nearest neighbor searching with applications to lattice sieving. In Robert Krauthgamer, editor, *27th SODA*, pages 10–24. ACM-SIAM, January 2016. 27
- [BDH11] Johannes A. Buchmann, Erik Dahmen, and Andreas Hülsing. XMSS - A practical forward secure signature scheme based on minimal security assumptions. In *PQCrypto*, pages 117–129, 2011. 19
- [BG14] Shi Bai and Steven D. Galbraith. An improved compression technique for signatures based on learning with errors. In *CT-RSA*, pages 28–47, 2014. 1, 2, 5, 13, 20, 27



- [BHH<sup>+</sup>15] Daniel J. Bernstein, Daira Hopwood, Andreas Hülsing, Tanja Lange, Ruben Niederhagen, Louiza Papachristodoulou, Michael Schneider, Peter Schwabe, and Zooko Wilcox-O’Hearn. SPHINCS: practical stateless hash-based signatures. In Marc Fischlin and Elisabeth Oswald, editors, *Advances in Cryptology – EUROCRYPT 2015*, volume 9056 of *LNCS*, pages 368–397. Springer, 2015. <http://cryptojedi.org/papers/#sphincs>. 21
- [BHL<sup>+</sup>16] Leon Groot Bruinderink, Andreas Hülsing, Tanja Lange, and Yuval Yarom. Flush, gauss, and reload - A cache attack on the BLISS lattice-based signature scheme. In *CHES*, pages 323–345, 2016. 1
- [BKW03] Avrim Blum, Adam Kalai, and Hal Wasserman. Noise-tolerant learning, the parity problem, and the statistical query model. *J. ACM*, 50(4):506–519, 2003. 30
- [BS16] Jean-François Biasse and Fang Song. Efficient quantum algorithms for computing class groups and solving the principal ideal problem in arbitrary degree number fields. In Robert Krauthgamer, editor, *27th SODA*, pages 893–902. ACM-SIAM, January 2016. 30
- [CDPR16] Ronald Cramer, Léo Ducas, Chris Peikert, and Oded Regev. Recovering short generators of principal ideals in cyclotomic rings. In Marc Fischlin and Jean-Sébastien Coron, editors, *EUROCRYPT 2016, Part II*, volume 9666 of *LNCS*, pages 559–585. Springer, Heidelberg, May 2016. 30
- [CDW17] Ronald Cramer, Léo Ducas, and Benjamin Wesolowski. Short Stickelberger class relations and application to ideal-SVP. In *EUROCRYPT (1)*, volume 10210 of *Lecture Notes in Computer Science*, pages 324–348, 2017. 30
- [CGS14] Peter Campbell, Michael Groves, and Dan Shepherd. Soliloquy: A cautionary tale. In *ETSI 2nd Quantum-Safe Crypto Workshop*, pages 1–9, 2014. 30
- [CHR<sup>+</sup>16] Ming-Shing Cheng, Andreas Hülsing, Joost Rijneveld, Simona Samardjiska, and Peter Schwabe. From 5-pass MQ-based identification to MQ-based signatures. In Jung Hee Cheon and Tsuyoshi Takagi, editors, *Advances in Cryptology – ASIACRYPT 2016*, volume 10032 of *LNCS*, page 135–165. Springer, 2016. <http://cryptojedi.org/papers/#mqdss>. 19, 21
- [CLP<sup>+</sup>17] Ming-Shing Chen, Wen-Ding Li, Bo-Yuan Peng, Bo-Yin Yang, and Chen-Mou Cheng. Implementing 128-bit secure mpkc signatures. Cryptology ePrint Archive, Report 2017/636, 2017. <http://eprint.iacr.org/2017/636/>. 21
- [CN11] Yuanmi Chen and Phong Q. Nguyen. BKZ 2.0: Better lattice security estimates. In Dong Hoon Lee and Xiaoyun Wang, editors, *ASIACRYPT 2011*, volume 7073 of *LNCS*, pages 1–20. Springer, Heidelberg, December 2011. 27
- [DDLL13] Léo Ducas, Alain Durmus, Tancrede Lepoint, and Vadim Lyubashevsky. Lattice signatures and bimodal gaussians. In *CRYPTO (1)*, pages 40–56, 2013. 1, 20, 21
- [DFSS07] Vivien Dubois, Pierre-Alain Fouque, Adi Shamir, and Jacques Stern. Practical cryptanalysis of SFLASH. In Alfred Menezes, editor, *Advances in Cryptology – CRYPTO 2007*, volume 4622 of *LNCS*, pages 1–12. Springer, 2007. <https://eprint.iacr.org/2007/14>. 19
- [DLP14] Léo Ducas, Vadim Lyubashevsky, and Thomas Prest. Efficient identity-based encryption over NTRU lattices. In *ASIACRYPT*, pages 22–41, 2014. 1, 20, 21

- [EFGT17] Thomas Espitau, Pierre-Alain Fouque, Benoit Gerard, and Mehdi Tibouchi. Side-channel attacks on bliss lattice-based signatures – exploiting branch tracing against strongswan and electromagnetic emanations in microcontrollers. Cryptology ePrint Archive, Report 2017/505, 2017. <http://eprint.iacr.org/2017/505> To appear in CCS 2017. 1
- [FGP<sup>+</sup>15] Jean-Charles Faugère, Danilo Gligoroski, Ludovic Perret, , Simona Samardjiska, and Enrico Thomae. A polynomial-time key-recovery attack on MQQ cryptosystems. In Jonathan Katz, editor, *Public-Key Cryptography – PKC 2015*, volume 9020 of *LNCS*, pages 150–174. Springer, 2015. <https://eprint.iacr.org/2014/811>. 19
- [GLP12] Tim Güneysu, Vadim Lyubashevsky, and Thomas Pöppelmann. Practical lattice-based cryptography: A signature scheme for embedded systems. In *CHES*, pages 530–547, 2012. 1, 2, 5, 20
- [GOPS13] Tim Güneysu, Tobias Oder, Thomas Pöppelmann, and Peter Schwabe. Software speed records for lattice-based signatures. In *PQCrypto*, volume 7932 of *Lecture Notes in Computer Science*, pages 67–82. Springer, 2013. 20
- [GPV08] Craig Gentry, Chris Peikert, and Vinod Vaikuntanathan. Trapdoors for hard lattices and new cryptographic constructions. In *STOC*, pages 197–206, 2008. 20
- [GQ88] Louis C. Guillou and Jean-Jacques Quisquater. A "paradoxical" indentity-based signature scheme resulting from zero-knowledge. In *CRYPTO*, pages 216–231, 1988. 11, 14
- [HHGP<sup>+</sup>03] Jeffrey Hoffstein, Nick Howgrave-Graham, Jill Pipher, Joseph H. Silverman, and William Whyte. Ntrusign: Digital signatures using the ntru lattice. In *CT-RSA*, pages 122–140, 2003. 20
- [HPS11] Guillaume Hanrot, Xavier Pujol, and Damien Stehlé. Analyzing blockwise lattice algorithms using dynamical systems. In Phillip Rogaway, editor, *CRYPTO 2011*, volume 6841 of *LNCS*, pages 447–464. Springer, Heidelberg, August 2011. 27
- [KF17] Paul Kirchner and Pierre-Alain Fouque. Revisiting lattice attacks on overstretched NTRU parameters. In *EUROCRYPT (1)*, volume 10210 of *Lecture Notes in Computer Science*, pages 3–26, 2017. 20, 30
- [KLS17] Eike Kiltz, Vadim Lyubashevsky, and Christian Schaffner. A concrete treatment of fiat-shamir signatures in the quantum random-oracle model. Cryptology ePrint Archive, Report 2017/916, 2017. <http://eprint.iacr.org/2017/916>. To appear in Eurocrypt 2018. 5, 12, 13, 27
- [KS98] Aviad Kipnis and Adi Shamir. Cryptanalysis of the Oil & Vinegar signature scheme. In Hugo Krawczyk, editor, *Advances in Cryptology – CRYPTO '98*, volume 1462 of *LNCS*, pages 257–266. Springer, 1998. 19
- [Laa15] Thijs Laarhoven. *Search problems in cryptography*. PhD thesis, Eindhoven University of Technology, 2015. 27
- [LN16] Patrick Longa and Michael Naehrig. Speeding up the number theoretic transform for faster ideal lattice-based cryptography. In *Cryptology and Network Security - 15th International Conference, CANS 2016, Milan, Italy, November 14-16, 2016, Proceedings*, pages 124–139, 2016. 4, 5

- [Lyu09] Vadim Lyubashevsky. Fiat-Shamir with aborts: Applications to lattice and factoring-based signatures. In *ASIACRYPT*, pages 598–616, 2009. 1, 2, 5, 13, 20
- [Lyu12] Vadim Lyubashevsky. Lattice signatures without trapdoors. In *EUROCRYPT*, pages 738–755, 2012. 5, 13, 20, 27
- [MW17] Daniele Micciancio and Michael Walter. Gaussian sampling over the integers: Efficient, generic, constant-time. In *CRYPTO*, pages 455–485, 2017. 20
- [PBY17] Peter Pessl, Leon Groot Bruinderink, and Yuval Yarom. To bliss-b or not to be - attacking strongswan’s implementation of post-quantum signatures. Cryptology ePrint Archive, Report 2017/490, 2017. <http://eprint.iacr.org/2017/490>. To appear in CCS 2017. 1
- [Sch89] Claus-Peter Schnorr. Efficient identification and signatures for smart cards. In *CRYPTO*, pages 239–252, 1989. 11, 14
- [SE94] Claus-Peter Schnorr and M. Euchner. Lattice basis reduction: Improved practical algorithms and solving subset sum problems. *Math. Program.*, 66:181–199, 1994. 27
- [Sei18] Gregor Seiler. Faster avx2 optimized ntt multiplication for ring-lwe lattice cryptography. Cryptology ePrint Archive, Report 2018/039, 2018. <https://eprint.iacr.org/2018/039>. 4, 5
- [TW12] Enrico Thomae and Christopher Wolf. Cryptanalysis of enhanced TTS, STS and all its variants, or: Why cross-terms are important. In Aikaterini Mitro-kotsa and Serge Vaudenay, editors, *Progress in Cryptology – AFRICACRYPT 2012*, volume 7374 of *LNCS*, pages 188–202. Springer, 2012. 19

## A Proofs for Rounding Algorithm Properties

The three lemmas below prove each of the three parts of Lemma 1.

**Lemma 3.** *Let  $r, z \in \mathbb{Z}_q$  with  $\|z\|_\infty \leq \alpha/2$ . Then*

$$\text{UseHint}_q(\text{MakeHint}_q(z, r, \alpha), r, \alpha) = \text{HighBits}_q(r + z, \alpha).$$

*Proof.* The output of  $\text{Decompose}_q$  is an integer  $r_1$  such that  $0 \leq r_1 < (q - 1)/\alpha$  and another integer  $r_0$  such that  $\|r_0\|_\infty \leq \alpha/2$ . Because  $\|z\|_\infty \leq \alpha/2$ , the integer  $v_1 := \text{HighBits}_q(r + z, \alpha)$  either stays the same as  $r_1$  or becomes  $r_1 \pm 1$  modulo  $m = (q - 1)/\alpha$ . More precisely, if  $r_0 > 0$ , then  $-\alpha/2 < r_0 + z \leq \alpha$ . This implies that  $v_1$  is either  $r_1$  or  $r_1 + 1 \bmod m$ . If  $r_0 \leq 0$ , then we have  $-\alpha \leq r_0 + z \leq \alpha/2$ . In this case, we have  $v_1 = r_1$  or  $r_1 - 1 \bmod m$ .

The  $\text{MakeHint}_q$  routine checks whether  $r_1 = v_1$  and outputs 0 if this is so, and 1 if  $r_1 \neq v_1$ . The  $\text{UseHint}_q$  routine uses the “hint”  $h$  to either output  $r_1$  (if  $y = 0$ ) or, depending on whether  $r_0 > 0$  or not, output either  $r_1 + 1 \bmod^+ m$  or  $r_1 - 1 \bmod^+ m$ .  $\square$

The lemma below shows that  $r$  is not too far away from the output of the  $\text{UseHint}_q$  algorithm. This will be necessary for the security of the scheme.

**Lemma 4.** *Let  $(h, r) \in \{0, 1\} \times \mathbb{Z}_q$  and let  $v_1 = \text{UseHint}_q(h, r, \alpha)$ . If  $h = 0$ , then  $\|r - v_1 \cdot \alpha\|_\infty \leq \alpha/2$ ; else  $\|r - v_1 \cdot \alpha\|_\infty \leq \alpha + 1$ .*

*Proof.* Let  $(r_1, r_0) := \text{Decompose}_q(r, \alpha)$ . We go through all three cases of the  $\text{UseHint}_q$  procedure.

Case 1 ( $h = 0$ ): We have  $v_1 = r_1$  and

$$r - v_1 \cdot \alpha = r_1 \cdot \alpha + r_0 - r_1 \cdot \alpha = r_0,$$

which by definition has absolute value at most  $\alpha/2$ .

Case 2 ( $h = 1$  and  $r_0 > 0$ ): We have  $v_1 = r_1 + 1 - \kappa \cdot (q-1)/\alpha$  for  $\kappa = 0$  or  $1$ . Thus

$$\begin{aligned} r - v_1 \cdot \alpha &= r_1 \cdot \alpha + r_0 - (r_1 + 1 - \kappa \cdot (q-1)/\alpha) \cdot \alpha \\ &= -\alpha + r_0 + \kappa \cdot (q-1). \end{aligned}$$

After centered reduction modulo  $q$ , the latter has magnitude  $\leq \alpha$ .

Case 3 ( $h = 1$  and  $r_0 \leq 0$ ): We have  $v_1 = r_1 - 1 + \kappa \cdot (q-1)/\alpha$  for  $\kappa = 0$  or  $1$ . Thus

$$\begin{aligned} r - v_1 \cdot \alpha &= r_1 \cdot \alpha + r_0 - (r_1 - 1 + \kappa \cdot (q-1)/\alpha) \cdot \alpha \\ &= \alpha + r_0 - \kappa \cdot (q-1). \end{aligned}$$

After centered reduction modulo  $q$ , the latter quantity has magnitude  $\leq \alpha + 1$ .  $\square$

The next lemma will play a role in proving the strong existential unforgeability of our signature scheme. It states that two different  $h, h'$  cannot lead to  $\text{UseHint}_q(h, r, \alpha) = \text{UseHint}_q(h', r, \alpha)$ .

**Lemma 5.** *Let  $r \in \mathbb{Z}_q$  and  $h, h' \in \{0, 1\}$ . If  $\text{UseHint}_q(h, r, \alpha) = \text{UseHint}_q(h', r, \alpha)$ , then  $h = h'$ .*

*Proof.* Note that  $\text{UseHint}_q(0, r, \alpha) = r_1$  and  $\text{UseHint}_q(1, r, \alpha)$  is equal to  $(r_1 \pm 1) \bmod^+(q-1)/\alpha$ . Since  $(q-1)/\alpha \geq 2$ , we have that  $r_1 \neq (r_1 \pm 1) \bmod^+(q-1)/\alpha$ .  $\square$

We now prove Lemma 2.

*Proof.* (Of Lemma 2) We prove the lemma for integers, rather than vectors of polynomials, since the  $\text{HighBits}$  function works independently on each coefficient. If  $\|\text{LowBits}_q(r, \alpha)\|_\infty < \alpha/2 - \beta$ , then  $r = r_1 \cdot \alpha + r_0$  where  $-\alpha/2 + \beta < r_0 \leq \alpha/2 + \beta$ . Then  $r + s = r_1 \cdot \alpha + (r_0 + s)$  and  $-\alpha/2 < r_0 + s \leq \alpha/2$ . Therefore  $r + s \bmod^\pm \alpha = r_0 + s$ , and thus

$$(r + s) - ((r + s) \bmod^\pm \alpha) = r_1 \cdot \alpha = r - (r \bmod^\pm \alpha),$$

and the claim in the Lemma follows.  $\square$

## B Zero-Knowledge Proof

The security of our scheme does not rely on the part of the public key  $\mathbf{t}_0$  being secret and so we will be assuming that the public key is  $\mathbf{t}$  rather than  $\mathbf{t}_1$ .

We want to first compute the probability that some particular  $(\mathbf{z}, c)$  is generated in Step 17 taken over the randomness of  $\mathbf{y}$  and the random oracle  $H$  which is modeled as a random function. We have

$$\Pr[\mathbf{z}, c] = \Pr[c] \cdot \Pr[\mathbf{y} = \mathbf{z} - c\mathbf{s}_1 \mid c].$$

Whenever  $\mathbf{z}$  has all its coefficients less than  $\gamma_1 - \beta$  then the above probability is exactly the same for every such tuple  $(\mathbf{z}, c)$ . This is because  $\|c\mathbf{s}_i\|_\infty \leq \beta$  (with overwhelming probability), and thus  $\|\mathbf{z} - c\mathbf{s}_1\|_\infty \leq \gamma_1 - 1$ , which is a valid value of  $\mathbf{y}$ . Therefore, if we

only output  $\mathbf{z}$  when all its coefficients have magnitudes less than  $\gamma_1 - \beta$ , then the resulting distribution will be uniformly random over  $S_{\gamma_1 - \beta - 1}^\ell \times B_{60}$ .

The simulation of the signature follows [Lyu12, BG14]. The simulator picks a uniformly random  $(\mathbf{z}, c)$  in  $S_{\gamma_1 - \beta - 1}^\ell \times B_{60}$ , after which it also makes sure that

$$\|\mathbf{r}_0\|_\infty = \|\text{LowBits}_q(\mathbf{w} - c\mathbf{s}_2, 2\gamma_2)\|_\infty < \gamma_2 - \beta.$$

By Equation (2), we know that  $\mathbf{w} - c\mathbf{s}_2 = \mathbf{A}\mathbf{z} - c\mathbf{t}$ , and therefore the simulator can perfectly simulate this as well.

If  $\mathbf{z}$  does indeed satisfy  $\|\text{LowBits}_q(\mathbf{w} - c\mathbf{s}_2, 2\gamma_2)\|_\infty < \gamma_2 - \beta$ , then as long as  $\|c\mathbf{s}_2\|_\infty \leq \beta$ , we will have

$$\mathbf{r}_1 = \text{HighBits}_q(\mathbf{w} - c\mathbf{s}_2, 2\gamma_2) = \text{HighBits}_q(\mathbf{w}, 2\gamma_2) = \mathbf{w}_1.$$

Since our  $\beta$  was chosen such that the probability (over the choice of  $c, \mathbf{s}_2$ ) that  $\|c\mathbf{s}_2\|_\infty < \beta$  is  $> 1 - 2^{-128}$ , the simulator does not need to perform the check that  $\mathbf{r}_1 = \mathbf{w}_1$  and can always assume that it passes.

We can then program

$$H(\mu \parallel \mathbf{w}_1) \leftarrow c.$$

Unless we have already set the value of  $H(\mu \parallel \mathbf{w}_1)$  to something else, the resulting pair  $(\mathbf{z}, c)$  has the same distribution as in a genuine signature of  $\mu$ . It was shown in [KLS17] that the probability, over the random choice of  $\mathbf{A}$  and  $\mathbf{y}$ , that we already set the value of  $H(\mu \parallel \mathbf{w}_1)$  is less than  $2^{-255}$ .

All the other steps (after Step 19) of the signing algorithm are performed using public information and are therefore simulatable.

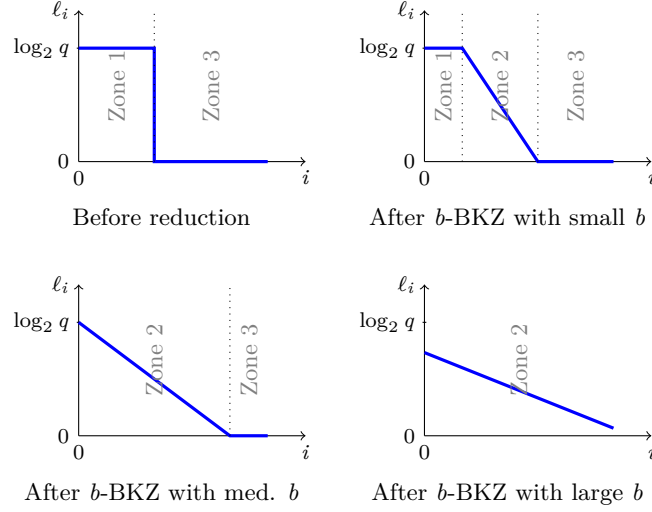
## C Concrete Security

### C.1 Lattice Reduction and Core-SVP Hardness

The best known algorithm for finding very short non-zero vectors in Euclidean lattices is the Block–Korkine–Zolotarev algorithm (BKZ) [SE94], proposed by Schnorr and Euchner in 1991. More recently, it was proven to quickly converge to its fix-point [HPS11] and improved in practice [CN11]. Yet, what it achieves asymptotically remains unchallenged.

BKZ with block-size  $b$  makes calls to an algorithm that solves the Shortest lattice Vector Problem (SVP) in dimension  $b$ . The security of our scheme relies on the necessity to run BKZ with a large block-size  $b$  and the fact that the cost of solving SVP is exponential in  $b$ . The best known classical SVP solver [BDGL16] runs in time  $\approx 2^{c_C \cdot b}$  with  $c_C = \log_2 \sqrt{3/2} \approx 0.292$ . The best known quantum SVP solver [Laa15, Sec. 14.2.10] runs in time  $\approx 2^{c_Q \cdot b}$  with  $c_Q = \log_2 \sqrt{13/9} \approx 0.265$ . One may hope to improve these run-times, but going below  $\approx 2^{c_P \cdot b}$  with  $c_P = \log_2 \sqrt{4/3} \approx 0.2075$  would require a theoretical breakthrough. Indeed, the best known SVP solvers rely on covering the  $b$ -dimensional sphere with cones of center-to-edge angle  $\pi/3$ : this requires  $2^{c_P \cdot b}$  cones. The subscripts C, Q, P respectively stand for Classical, Quantum and Paranoid.

The strength of BKZ increases with  $b$ . More concretely, given as input a basis  $(\mathbf{c}_1, \dots, \mathbf{c}_n)$  of an  $n$ -dimensional lattice, BKZ repeatedly uses the  $b$ -dimensional SVP-solver on lattices of the form  $(\mathbf{c}_{i+1}(i), \dots, \mathbf{c}_j(i))$  where  $i \leq n$ ,  $j = \min(n, i + b)$  and where  $\mathbf{c}_k(i)$  denotes the projection of  $\mathbf{c}_k$  orthogonally to the vectors  $(\mathbf{c}_1, \dots, \mathbf{c}_i)$ . The effect of these calls is to flatten the curve of the  $\ell_i = \log_2 \|\mathbf{c}_i(i-1)\|$ 's (for  $i = 1, \dots, n$ ). At the start of the execution, the  $\ell_i$ 's typically decrease fast, at least locally. As BKZ preserves the determinant of the  $\mathbf{c}_i$ 's, the sum of the  $\ell_i$ 's remains constant throughout the execution, and after a (small) polynomial number of SVP calls, BKZ has made the  $\ell_i$ 's decrease less.



**Figure 5:** Evolution of Gram-Schmidt length in log-scale under BKZ reduction for various block sizes. The area under the curves remains constant, and the slope in Zone 2 decrease with the blocksize. Note that Zone 3 may disappear before Zone 1, depending on the shape of the input basis.

It can be heuristically estimated that for sufficiently large  $b$ , the local slope of the  $\ell_i$ 's converges to

$$\text{slope}(b) = \frac{1}{b-1} \log_2 \left( \frac{b}{2\pi e} (\pi \cdot b)^{1/b} \right),$$

unless the local input  $\ell_i$ 's are already too small or too large. The quantity  $\text{slope}(b)$  decreases with  $b$ , implying that the larger  $b$  the flatter the output  $\ell_i$ 's.

In our case, the input  $\ell_i$ 's are of the following form (cf. Fig. 5). The first ones are all equal to  $\log_2 q$  and the last ones are all equal to 0. BKZ will flatten the jump, decreasing  $\ell_i$ 's with small  $i$ 's and increasing  $\ell_i$ 's with large  $i$ 's. However, the local slope  $\text{slope}(b)$  may not be sufficiently small to make the very first  $\ell_i$ 's decrease and the very last  $\ell_i$ 's increase. Indeed, BKZ will not increase (resp. increase) some  $\ell_i$ 's if these are already smaller (resp. larger) than ensured by the local slope guarantee. In our case, the  $\ell_i$ 's are always of the following form at the end of the execution:

- The first  $\ell_i$ 's are constant equal to  $\log_2 q$  (this is the possibly empty Zone 1).
- Then they decrease linearly, with slope  $\text{slope}(b)$  (this is the never-empty Zone 2).
- The last  $\ell_i$ 's are constant equal to 0 (this is the possibly empty Zone 3).

The graph is continuous, i.e., if Zone 1 (resp. Zone 3) is not empty, then Zone 2 starts with  $\ell_i = \log_2 q$  (resp. ends with  $\ell_i = 0$ ).

## C.2 Solving MLWE

Any MLWE $_{\ell,k,D}$  instance for some distribution  $D$  can be viewed as an LWE instance of dimensions  $256 \cdot \ell$  and  $256 \cdot k$ . Indeed, the above can be rewritten as finding  $\text{vec}(\mathbf{s}_1), \text{vec}(\mathbf{s}_2) \in \mathbb{Z}^{256 \cdot \ell} \times \mathbb{Z}^{256 \cdot k}$  from  $(\text{rot}(\mathbf{A}), \text{vec}(\mathbf{t}))$ , where  $\text{vec}(\cdot)$  maps a vector of ring elements to the vector obtained by concatenating the coefficients of its coordinates, and  $\text{rot}(\mathbf{A}) \in \mathbb{Z}_q^{256 \cdot k \times 256 \cdot \ell}$  is obtained by replacing all entries  $a_{ij} \in R_q$  of  $\mathbf{A}$  by the  $256 \times 256$  matrix whose  $z$ -th column is  $\text{vec}(x^{z-1} \cdot a_{ij})$ .

Given an LWE instance, there are two lattice-based attacks. The primal attack and the dual attack. Here, the primal attack consists in finding a short non-zero vector in the lattice  $\Lambda = \{\mathbf{x} \in \mathbb{Z}^d : \mathbf{M}\mathbf{x} = \mathbf{0} \bmod \mathbf{q}\}$  where  $\mathbf{M} = (\text{rot}(\mathbf{A})_{[1:m]} | \mathbf{I}_m | \text{vec}(\mathbf{t})_{[1:m]})$  is an  $m \times d$  matrix where  $d = 256 \cdot \ell + m + 1$  and  $m \leq 256 \cdot k$ . Indeed, it is sometime not optimal to use all the given equations in lattice attacks.

We tried all possible number  $m$  of rows, and, for each trial, we increased the blocksize of  $b$  until the value  $\ell_{d-b}$  obtained as explained above was deemed sufficiently large. As explained in [ADPS16, Sec. 6.3], if  $2^{\ell_{d-b}}$  is greater than the expected norm of  $(\text{vec}(\mathbf{s}_1), \text{vec}(\mathbf{s}_2))$  after projection orthogonally to the first  $d - b$  vectors, it is likely that the Module-LWE solution can be easily extracted from the BKZ output.

The dual attack consists in finding a short non-zero vector in the lattice  $\Lambda' = \{(\mathbf{x}, \mathbf{y}) \in \mathbb{Z}^m \times \mathbb{Z}^d : \mathbf{M}^T \mathbf{x} + \mathbf{y} = \mathbf{0} \bmod q\}$ ,  $\mathbf{M} = (\text{rot}(\mathbf{A})_{[1:m]})$  is an  $m \times d$  matrix where  $d = 256 \cdot \ell$  and  $m \leq 256 \cdot k$ . Again, for each value of  $m$ , we increased the value of  $b$  until the value  $\ell_1$  obtained as explained above was deemed sufficiently small according the analysis of [ADPS16, Sec. 6.3].

### C.3 Solving MSIS and SelfTargetMSIS

As per the discussion in Section 4.2.1, the best known attack against the SelfTargetMSIS problem involves either breaking the security of H or solving the problem in Eq. (12). The latter amounts to solving the MSIS $_{k,\ell+1,\zeta}$  problem for the matrix  $[\mathbf{A} | \mathbf{t}']$ .<sup>9</sup>

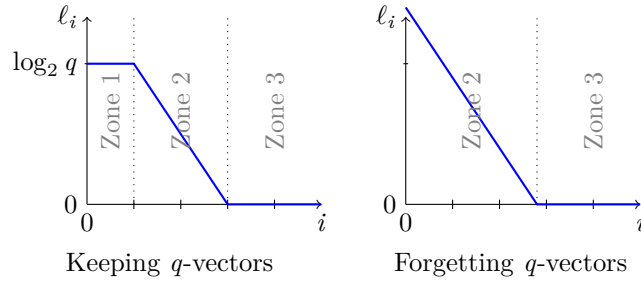
Note that the MSIS instance can be mapped to a SIS $_{256 \cdot k, 256 \cdot (\ell+1), \zeta}$  instance by considering the matrix  $\text{rot}(\mathbf{A} | \mathbf{t}') \in \mathbb{Z}_q^{256 \cdot k \times 256 \cdot (\ell+1)}$ . The attack against the MSIS $_{k,\ell,\zeta'}$  instance in Eq. (6) can similarly be mapped to a SIS $_{256 \cdot k, 256 \cdot \ell, \zeta'}$  instance by considering the matrix  $\text{rot}(\mathbf{A}) \in \mathbb{Z}_q^{256 \cdot k \times 256 \cdot \ell}$ . The attacker may consider a subset of  $w$  columns, and let the solution coefficients corresponding to the dismissed columns be zero.

*Remark 1.* An unusual aspect here is that we are considering the infinity norm, rather than the Euclidean norm. Further, for our specific parameters, the Euclidean norms of the solutions are above  $q$ . In particular, the vector  $(q, 0, \dots, 0)^T$  belongs to the lattice, has Euclidean norm below that of the solution, but its infinity norm above the requirement. This raises difficulties in analyzing the strength of BKZ towards solving our infinity norm SIS instances: indeed, even with small values of  $b$ , the first  $\ell_i$ 's are short (they correspond to  $q$ -vectors), even though they are not solutions.

For each number  $w$  of selected columns and for each value of  $b$ , we compute the estimated BKZ output  $\ell_i$ 's, as explained above. We then consider the smallest  $i$  such that  $\ell_i$  is below  $\log_2 q$  and the largest  $j$  such that  $\ell_j$  above 0. These correspond to the vectors that were modified by BKZ, with smallest and largest indices, respectively. In fact, for the same cost as a call to the SVP-solver, we can obtain  $\sqrt{4/3}^b$  vectors with Euclidean norm  $\approx 2^{\ell_i}$  after projection orthogonally to the first  $i - 1$  basis vectors. Now, let us look closely at the shape of such a vector. As the first  $i - 1$  basis vectors are the first  $i - 1$  canonical unit vectors multiplied by  $q$ , projecting orthogonally to these consists in zeroing the first  $i - 1$  coordinates. The remaining  $w - i + 1$  coordinates have total Euclidean norm  $\approx 2^{\ell_i} \approx q$ , and the last  $w - j$  coordinates are 0. We heuristically assume that these coordinates have similar magnitudes  $\sigma \approx 2^{\ell_i} / \sqrt{j - i + 1}$ ; we model each such coordinate as a Gaussian of standard deviation  $\sigma$ . We assume that each one of our  $\sqrt{4/3}^b$  vectors has its first  $i - 1$  coordinates independently uniformly distributed modulo  $q$ , and finally compute the probability that all coordinates in both ranges  $[0, i - 1]$  and  $[i, j]$  are less than  $B$  in absolute value. Our cost estimate is the inverse of that probability multiplied by the run-time of our  $b$ -dimensional SVP-solver.

<sup>9</sup>Note that a solution to Eq. (12) would require the coefficient in front of  $\mathbf{t}'$  to be  $\pm 1$ , while we're allowing any small polynomial. Furthermore, as discussed after Eq. (12), some parts of the real solution are smaller than the bound  $\zeta$ , but we're ignoring this for the sake of being conservative with our analysis.





**Figure 6:** Effect of forgetting  $q$ -vectors by randomization, under the same BKZ-blocksize  $b$ .

**Forgetting  $q$ -vectors.** For all the parameter sets proposed in this paper, the best parametrization of the attack above kept the basis in a shape with a non-trivial Zone 1. We note that the coordinates in this range have a quite lower probability of passing the  $\ell_\infty$  constraint than coordinates in Zone 2. We therefore considered a strategy consisting of “forgetting” the  $q$ -vectors, by re-randomizing the input basis before running the BKZ algorithm. For the same blocksize  $b$ , this makes Zone 1 of the output basis disappear (BKZ does not find the  $q$ -vectors), at the cost of producing a basis with first vectors of larger Euclidean norms. This is depicted in Fig. 6.

It turns out that this strategy always improves over the previous strategy for the parameter ranges considered in this paper. We therefore used this strategy for our security estimates.

#### C.4 On Other Attacks

For our parameters, the BKW [BKW03] and Arora–Ge [AG11] families of algorithms are far from competitive.

**Algebraic attacks.** One specificity of our LWE and SIS instances is that they are inherited from Module-LWE and Module-SIS instances. One may wonder whether the extra algebraic structure of the resulting lattices can be exploited by an attacker. The line of work of [CGS14, BS16, CDPR16, CDW17] did indeed find new cryptanalytic results on certain algebraic lattices, but [CDW17] mentions serious obstacles towards breaking cryptographic instances of Ring-LWE. By switching from Ring-LWE to Module-LWE, we get even further away from those weak algebraic lattice problems.

**Dense sublattice attacks.** Kirchner and Fouque [KF17] showed that the existence of many linearly independent and unexpectedly short lattice vectors (much shorter than Minkowski’s bound) helps BKZ run better than expected in some cases. This could happen for our primal LWE attack, by extending  $\mathbf{M} = (\text{rot}(\mathbf{A})_{[1:m]} | \mathbf{I}_m | \text{vec}(\mathbf{t})_{[1:m]})$  to  $(\text{rot}(\mathbf{A})_{[1:m]} | \mathbf{I}_m | \text{rot}(\mathbf{t})_{[1:m]})$ : the associated lattice now has 256 linearly independent short vectors rather than a single one. The Kirchner-Fouque analysis of BKZ works best if both  $q$  and the ratio between the number of unexpectedly short vectors and the lattice dimension are high. In the NTRU case, for example, the ratio is  $1/2$ , and, for some schemes derived from NTRU, the modulus  $q$  is also large. We considered this refined analysis of BKZ in our setup, but, to become relevant for our parameters, it requires a parameter  $b$  which is higher than needed with the usual analysis of BKZ. Note that [KF17] also arrived to the conclusion that this attack is irrelevant in the small modulus regime, and is mostly a threat to fully homomorphic encryption schemes and cryptographic multilinear maps.

Note that, once again, the switch from Ring-LWE to MLWE takes us further away from lattices admitting unconventional attacks. Indeed, the dimension ratio of the dense

sub-lattice is  $1/2$  in NTRU, at most  $1/3$  in lattices derived from Ring-LWE, and at most  $1/(\ell + 2)$  in lattices derived from MLWE.

**Specialized attack against  $\ell_\infty$ -SIS.** At last, we would like to mention that it is not clear whether the attack sketched in Appendix C.3 above for SIS in infinity norm is optimal. Indeed, as we have seen, this approach produces many vectors, with some rather large uniform coordinates (at indices  $1, \dots, i$ ), and smaller Gaussian ones (at indices  $i, \dots, j$ ). In our current analysis, we simply hope that one of the vector satisfies the  $\ell_\infty$  bound. Instead, one could combine them in ways that decrease the size of the first (large) coefficients, while letting the other (small) coefficients grow a little bit.

This situation created by the use of  $\ell_\infty$ -SIS (see Remark 1) has — to the best of our knowledge — not been studied in detail. After a preliminary analysis, we do not consider such an improved attack a serious threat to our concrete security claims, especially in the light of the approximations already made in the favor of the adversary. Nevertheless, we believe this question deserves a detailed study, which we leave to future work.