

Cryptanalysis of the Revised NTRU Signature Scheme

Craig Gentry¹ and Mike Szydlo²

¹ DoCoMo USA Labs, San Jose, CA, USA
cgentry@docomolabs-usa.com

² RSA Laboratories, Bedford, MA, USA
mszydlo@rsasecurity.com

Keywords: NSS, NTRU, NTRUSign, Signature Scheme, Lattice Reduction, Cryptanalysis, Orthogonal Lattice, Cyclotomic Integer, Galois Congruence.

Abstract. In this paper, we describe a three-stage attack against Revised NSS, an NTRU-based signature scheme proposed at the Eurocrypt 2001 conference as an enhancement of the (broken) proceedings version of the scheme. The first stage, which typically uses a transcript of only 4 signatures, effectively cuts the key length in half while completely avoiding the intended hard lattice problem. After an empirically fast second stage, the third stage of the attack combines lattice-based and congruence-based methods in a novel way to recover the private key in polynomial time. This cryptanalysis shows that a passive adversary observing only a few valid signatures can recover the signer's entire private key. We also briefly address the security of NTRUSign, another NTRU-based signature scheme that was recently proposed at the rump session of Asiacrypt 2001. As we explain, some of our attacks on Revised NSS may be extended to NTRUSign, but a *much* longer transcript is necessary. We also indicate how the security of NTRUSign is based on the hardness of several problems, not solely on the hardness of the usual NTRU lattice problem.

1 Introduction

The Revised NTRU Signature Scheme (R-NSS) and “NTRUSign” are the two most recent of several signature schemes related to the NTRU encryption scheme (now called NTRUEncrypt). NTRUEncrypt and the related signature schemes are not based on traditional hard problems such as factoring or computing discrete logarithms, like much of today's cryptography. Instead, NTRUEncrypt was originally conceived as a cryptosystem based on polynomial arithmetic. Based on an early attack found by Coppersmith and Shamir [7], however, the underlying hard problem was soon reformulated as a lattice problem. See [22] for an update on how lattices have recently been used both as a cryptanalytic tool and as a potential basis for cryptography.

There are two reasons for seeking alternative hard problems on which cryptography may be based. First, it is prudent to hedge against the risk of potential

breakthroughs in factoring and computing discrete logarithms. A second and more significant reason is efficiency. NTRU-based algorithms, for example, are touted to run hundreds of times faster while providing the same security as competing algorithms. The drawback in using alternative hard problems is that they may not be as well understood. Although lattice theory has been studied for over 100 years,¹ the *algorithmic* nature of hard lattice problems such as the “shortest vector problem” (SVP) was not really studied *intensively* until Lenstra, Lenstra and Lovász discovered a polynomial-time lattice basis reduction algorithm in 1982. Moreover, NTRU-based schemes use specific types of lattices based on an underlying polynomial ring, and these lattices generate specific types of lattice problems that may be easier to solve than general lattice problems. Since these specific lattice problems have been studied intensively only since NTRUEncrypt’s introduction in 1996, we can expect plenty of new results. This paper is a case in point: we use a new polynomial-time algorithm to find the shortest vector in certain lattices that arise in R-NSS, allowing us to break the scheme.

1.1 History of NTRU-based Signature Schemes

Since the invention of NTRUEncrypt in 1996, several related identification and signature schemes have been proposed. In 1997, Hoffstein, Kaliski, et.al. filed a patent for an identification scheme based on constrained polynomials [12], although the scheme was later determined not to be secure. A “preliminary” version of NSS was first presented at the rump session of Crypto 2000, but this scheme was observed by its authors to be insecure at an early stage and by Mironov [21] independently a few months later. Essentially, the signatures leaked information about the private key, which was revealed by a straightforward averaging attack.

Certain adaptations were made to eliminate the disclosed weaknesses, yielding the scheme described in the proceedings of Eurocrypt 2001 [17]. (See also [16] and [4].) Even so, these signatures still leaked information about the private key: specifically, correlations between certain coefficients in the signature and the private key were sufficient to recover the entire public key. Moreover, the scheme was susceptible to a simple direct forgery attack through which an attacker could quickly sign arbitrary messages without even having to see a single legitimate signature. These forgery and key recovery attacks were presented by Gentry, Jonsson, Stern and Szydlo at the rump session of Eurocrypt 2001, the same conference where NSS was to be first fully presented, and were further described in an Asiacrypt 2001 paper [9]. They render the scheme as presented in [17], [16] and [4] completely insecure.

Informed of the attacks prior to the conference, the authors of NSS sketched a revised scheme in their Eurocrypt presentation. They described these revisions in more detail in a technical note entitled “Enhanced Encoding and Verification Methods for the NTRU Signature Scheme” [14], which was revised several months later [15]. They finally committed to a scheme, which we will call

¹ This includes early work by Hermite and Minkowski, the latter calling the topic “Geometrie der Zahlen” (Geometry of Numbers) in 1910.

“R-NSS,” by publishing it in the preliminary cryptographic standard document EESS [5]. They also published analysis and research showing how the new scheme defeated previous attacks. Although R-NSS does indeed appear to be a significantly stronger scheme than previous versions, this paper describes how it can be broken.

Since the initial submission of this paper, NTRU has proposed a new NTRU-based signature scheme called NTRUSign. Although our primary focus is R-NSS, we also provide security analysis of NTRUSign, as requested by the Program Committee.

1.2 Our Cryptanalysis

In our cryptanalysis of R-NSS, we use for concreteness the parameters suggested in the technical note [15] and standards document [5]. We show how a passive adversary who observes only a few valid signatures can recover the signer’s entire private key. Although some might consider R-NSS to be even more *ad hoc* than previous NTRU-based signature schemes, our attacks against it are more fundamental than previous attacks, in that they target the basic tenets of the scheme rather than its peculiarities.

The rest of this paper is organized as follows: In Section 2, we provide background mathematics, and then in Section 3, we describe R-NSS. In Section 4, we survey the previous attacks on NTRU-based signature schemes that are relevant to our cryptanalysis of R-NSS. In Section 5, we detail the first stage of our attack: the lifting procedure. Next, in Section 6, we describe how to obtain the polynomial $f * \bar{f}$, which we use in the final stage of the attack. In Section 7, we introduce novel techniques for circulant lattices which enable a surprising algorithm to obtain the private key f in polynomial time. We give a summary of our R-NSS cryptanalysis in Section 8. Finally, in Section 9, we describe NTRUSign, consider attacks against it and describe alternative hard problems that underlie its security.

2 Background Mathematics

As with NTRUEncrypt and previous NTRU-based signature schemes, the key underlying structure of R-NSS is the polynomial ring

$$R = \mathbb{Z}[X]/(X^N - 1) \tag{1}$$

where N is a prime integer (e.g., 251) in practice. In some steps, R-NSS uses the quotient ring $R_q = \mathbb{Z}_q[X]/(X^N - 1)$, where the coefficients are reduced modulo q , and normally taken in the range $(-q/2, q/2]$, where q is typically a power of 2 (e.g., 128).

Multiplication in R is similar to ordinary polynomial multiplication, but subject to the relations $X^{N+k} = X^k$ for any $k \geq 0$. This means that the coefficient of X^k in the product $a * b$ of $a = a_0 + a_1X + \dots + a_{N-1}X^{N-1}$ and

$b = b_0 + b_1X + \dots + b_{N-1}X^{N-1}$ is

$$(a * b)_k = \sum_{i+j=k \bmod N} a_i b_j. \quad (2)$$

The multiplication of two polynomials in R is also called the *convolution product* of the two polynomials. To any polynomial $a \in R$, it is also convenient to associate a *convolution matrix* to a as follows: Let M_a be the $N \times N$ circulant matrix indexed by $\{0, \dots, N-1\}$, where the element on position (i, j) is equal to $a_{(j-i) \bmod N}$. The product of a and b can also be expressed as the product of the row vector (a_0, \dots, a_{N-1}) with the matrix M_b . Any polynomial in R or R_q can be naturally associated with a row vector in this way, and we make this identification throughout the paper. We also use this identification to define the *Euclidean norm* of a polynomial:

$$\|f\| = \sqrt{\sum f_i^2}. \quad (3)$$

At times, we will refer to the reversal \bar{a} of a polynomial a , defined by $\bar{a}_k = a_{N-k}$ (with $\bar{a}_0 = a_0$). The mapping $a \mapsto \bar{a}$ is an automorphism of R , since applying the map twice yields the original polynomial. We use the term “palindromes” in referring to polynomials that are fixed under the reversal mapping on R – i.e., polynomials a such that $a = \bar{a}$. For any $a \in R$, it is easy to see that the product $a * \bar{a}$ is a palindrome. This fact, as well as the reversal mapping, may be described in elementary terms, but also in terms of an automorphism of the underlying cyclotomic field $\mathbb{Q}(\zeta_N)$. We refer the reader to Appendix D for further details on the Galois theory of R and $\mathbb{Q}(\zeta_N)$.

2.1 Lattices

The analysis of R-NSS will make frequent use of *lattices*. Formally, a lattice is a discrete subgroup of a vector space, but concretely, a lattice may be presented as the integral span of some set $B = \{b_0, \dots, b_{m-1}\}$ of linearly independent vectors in \mathbb{R}^N – that is,

$$L = \{v \mid v = \sum a_i b_i \mid a_i \in \mathbb{Z}\}. \quad (4)$$

We call m the dimension of the lattice, and B a basis of L . Bases will often be presented as a matrix in which the rows are the basis vectors $\{b_i\}$. Each lattice has an infinite number of bases, related by $B' = UB$ where U is a unimodular matrix, but some bases are more useful than others. The goal of *lattice reduction* is to find useful bases, typically ones with reasonably short, reasonably orthogonal basis vectors. The most celebrated lattice reduction algorithm is LLL [19], which has found many uses in cryptology. The contemporary survey [22] provides an overview of lattice techniques and [2] provides detailed descriptions of LLL variants.

The most famous lattice problem is the shortest vector problem (SVP): given a basis of a lattice L , find the shortest nonzero vector in L . Although LLL

and its variants manage to find *somewhat short* vectors in lattices, they do not necessarily find the *shortest* vector. In fact, SVP is an NP-hard problem (under randomized reductions) [1]. In previous cryptanalysis of NTRU and NSS, LLL’s inability to recover the shortest (or a very, very short) vector was a significant shortcoming. In some of our attacks, however, we will construct lattices in which even vectors that are only somewhat short reveal information about the signer’s private key, and then we will use LLL and its variants as black box algorithms to find these vectors. We will explain other aspects of lattice theory as they become relevant.

2.2 Ideals

Since R-NSS operates with polynomials in the ring R , we will need to consider multiplication in R , as well as ideals in this ring. Recall that an *ideal* is an additive subgroup of a commutative ring which is also closed under multiplication by any element in R , and a *principal ideal* is an ideal of R consisting of all R -multiples of a single element. We write (a) to denote the principal ideal consisting of all R -multiples of a , and say that the ideal is generated by a . We remark that not all ideals are principal, and furthermore, a generator of a principal ideal is not unique since there are infinitely many units $u \in R$, and for each u , both M_f and $M_{(f*u)}$ define the same ideal. We naturally extend these notions to lattices by defining a *lattice ideal* to be a lattice which is also closed under “multiplication” by polynomials in R , and a *principal ideal* to be a lattice which consists of all R -multiples of a given polynomial. Consider, for example, the lattice $L(M_a)$ generated by the circulant matrix M_a . Since the rows of M_a correspond to the various “rotations” $a * X^K$ of the element a , every possible a -multiple in R is in the lattice $L(M_a)$. Indeed, $L(M_a)$ is precisely the lattice equivalent of the ideal (a) . In this paper, we refer to these special lattices and ideals interchangeably.

In general, lattices are only closed under addition, and if multiplication is even defined, there is no guarantee that the product of two vectors will be another element of R that is also in the lattice. However, lattices corresponding to ideals in R do have this property. The algebraic structure of ideals is richer than that of general lattices, and we will exploit this extra structure in our novel attacks on R-NSS.

3 Description of R-NSS

The signature scheme R-NSS is a triplet (*keygen*, *sign*, *verify*) of algorithms operating on polynomials in $R = \mathbb{Z}[X]/(X^N - 1)$ and $R_q = \mathbb{Z}_q[X]/(X^N - 1)$, where N is prime, and $q < N$, (e.g. $N = 251$, $q = 128$). Other parameters in R-NSS include the modulus p , which is relatively prime to q and is typically chosen to be 3, as well as the integers d_u , d_f , d_g , d_m and d_z , whose suggested values are respectively 88, 52, 36, 80, and 58. The latter parameters are used to define several families of *ternary* polynomials as follows: $\mathcal{L}(d_1, d_2)$ denotes the set of polynomials in R_q , with d_1 coefficients 1, d_2 coefficients -1 and all other coefficients 0.

Key generation: Two polynomials f and g are randomly generated according to the equations

$$f = u + pf_1$$

$$g = u + pg_1$$

where $u \in \mathcal{L}(d_u, d_u + 1)$, $f_1 \in \mathcal{L}(d_f, d_f)$ and $g_1 \in \mathcal{L}(d_g, d_g)$. The signer keeps these two polynomials secret, with f serving as the signer's private key. The public key h is computed as $f^{-1} * g$ in R_q , and it is therefore necessary that f be invertible in R_q (i.e., $f * f^{-1} = 1$ for some $f^{-1} \in R_q$). This is true with very high probability (see [24]); in any case the preceding step may be repeated by choosing a different polynomial f_1 .

As in previous versions of NSS, the coefficients of f and g are small – i.e., they lie in a narrow range $([-4, 4])$ assuming $p = 3$) of \mathbb{Z}_q . However, R-NSS introduces a new secret polynomial – namely, u – into the private key generation process. In the previous version of NSS, $f \pmod{p}$ and $g \pmod{p}$ were public, allowing a statistical attack on a transcript of signatures. In Appendix B, we briefly describe this transcript attack, and explain how using u defeats it.

Signature generation: To sign a message, one transforms the message to be signed into a message representative according to a hash function-based procedure such as that described in [4]. We do not base any attack on this encoding, which can be made as safe as for any signature scheme. This message representative m is a polynomial in $\mathcal{L}(d_m, d_m)$. The signer then computes the following temporary variables:

$$y = u^{-1} * m \pmod{p}$$

$$z \in \mathcal{L}(d_z, d_z)$$

$$w = y + pz ,$$

where u^{-1} is computed in $R_p = \mathbb{Z}_p[X]/(X^N - 1)$. Notice that in R , $f * w \equiv m \pmod{p}$. This is not necessarily the case in R_q , however, since reduction modulo q causes “deviations” in the modulo p congruence, given that p and q are relatively prime. During the rest of the signing process, the signer will try to keep the number of these deviations to a minimum. The signer computes more temporary variables:

$$s = f * w \pmod{q}$$

$$t = g * w \pmod{q}$$

$$Dev_s = (s - m) \pmod{p}$$

$$Dev_t = (t - m) \pmod{p} .$$

Dev_s and Dev_t represent the deviations that the signer would like to correct, but, unfortunately for the signer, correcting a coefficient in s may cause additional deviations in t . Therefore, the signer limits his corrections to coefficient positions j such that $(Dev_s)_j = (Dev_t)_j$. He initializes a polynomial e to 0 and sets e_j to $-(Dev_s)_j$ when $(Dev_s)_j = (Dev_t)_j$. He then lets $e' = u^{-1} * e \pmod{p}$, adds e'

to w , and recomputes $s = f * w$ in R_q . The pair (m, s) is the signer’s signature of m .

The signing procedure above is not very intuitive, so we will pause at this point and try to explain the motivation behind it. The signer knows a pair of polynomials (f, g) satisfying $f * h = g$ in R_q where the coefficients of f and g lie in a very narrow range (e.g., $[-4, 4]$) of $(-q/2, q/2]$. An attacker can easily generate polynomials (a, b) satisfying $a * h = b$ in R_q , but the L2-norm of $(a||b)$ (concatenated) will likely be much more than the L2-norm of $(f||g)$. The point of the signing procedure is to give the signer a way of proving that he knows a *short* solution (a, b) to $a * h = b$ in R_q . Here is the essential approach: the signer (or forger) must produce a pair of polynomials (s, t) with coefficients in $(-q/2, q/2]$ that simultaneously satisfies $s * h = t \pmod{q}$ as well as $s \approx m \pmod{p}$ and $t \approx m \pmod{p}$ where “ \approx ” means that there are a small number of “deviations.” For the legitimate signer, this is easy: he can, for example, set $w \in R_p$ (coefficients in $\{-1, 0, 1\}$) to be $u^{-1} * m$ and let $s = f * w \pmod{q}$ and $t = g * w \pmod{q}$. Since f, g and w each have short L2-norm, $f * w$ and $g * w$ will also have somewhat short L2-norm, meaning that only a few coefficients of $f * w$ and $g * w$ will fall outside of $(-q/2, q/2]$, and therefore only a few deviations in the mod- p congruences will occur. An attacker who does not possess a short (a, b) satisfying $a * h = b$ in R_q will have trouble minimizing the number of mod- p deviations.

Signature verification: To avoid the forgery attacks presented in [9], verification has become a rather complicated process involving up to 20 distinct steps, detailed in [6] and [15], which fall into three broad categories: the *Quartile Distribution* tests, the *Mod 3 Distribution* tests, and the *L2 Norm* tests. In essence, the verifier checks, respectively, that

1. the coefficients of s and $t = s * h \pmod{q}$ have a roughly normal distribution;
2. the coefficients of s and t deviating from m are few and have a certain distribution; and
3. the L2 norms of $s' = p^{-1}(s - m) \pmod{q}$, $t' = p^{-1}(t - m) \pmod{q}$ and $(s' || t')$ (concatenated) are below certain thresholds.

Since we do not focus on forgery attacks in this paper, we defer the details of the verification process to Appendix A. On average, the signer has to run the signing process two or three times to produce a valid signature.

4 Previous Attacks on NSS

In this section, we review some relevant known attacks against NTRUEncrypt and previous NTRU-based signature schemes. This review will help us explain our cryptanalysis of R-NSS, which will occasionally leverage pieces of the attacks mentioned here.

4.1 Coppersmith-Shamir Attack

As with NTRUEncrypt, the security of R-NSS is claimed to be based on a hard lattice problem. Coppersmith and Shamir [7] were the first to present a lattice-based attack against NTRUEncrypt, an attack which is also relevant to R-NSS. Let L_{CS} be the lattice generated by the rows of the following matrix:

$$B_{CS} = \begin{bmatrix} I_{(N)} & M_h \\ 0 & qI_{(N)} \end{bmatrix}, \quad (5)$$

where $I_{(N)}$ is the N -dimensional identity matrix. This lattice clearly contains the vector $(f\|g)$, since $f * h = g \pmod{q}$. Moreover, for technical reasons [7], it is highly probable that $(f\|g)$ is the shortest nonzero vector in this lattice (up to rotation, sign, and excluding trivial vectors such as $(1^N, 1^N)$, the vector of all 1's.) Therefore, recovering the private key is simply a matter of recovering the shortest vector in L_{CS} , which we can presumably do using a lattice reduction algorithm. This attack is very effective when N is small (e.g., 107).

The problem with this approach (and lattice attacks, in general) is that no known lattice reduction algorithm is both very fast and very effective. More specifically, the LLL algorithm is “polynomial-time” – i.e., it terminates in time polynomial in the dimension m of the lattice – but, for high-dimensional lattices, such as those used in NTRU-based schemes, it almost certainly will *not* find the shortest vector. Rather, LLL only guarantees finding a vector that is no more than $2^{(m-1)/2}$ times as long as the shortest vector. Even though, in practice, LLL performs significantly better than this worst case bound, its performance is not sufficient for this lattice; we need the shortest vector or a *very* small multiple thereof. Other lattice reduction algorithms can find shorter vectors, but they naturally have greater time-complexity. The bottom line, based on current knowledge and on extensive empirical tests run by NTRU [23], seems to be that the time necessary to find $(f\|g)$ in L_{CS} grows at least exponentially in the dimension of the lattice ($2N$). This apparently hard lattice problem of recovering $(f\|g)$ from the $2N$ -dimensional lattice is claimed to underlie the security of both NTRUEncrypt and R-NSS.

Remark 1. As mentioned previously, the fact that SVP is an NP-hard problem for general lattices does not necessarily mean that finding short vectors in L_{CS} is hard. May [20] exploited the specifics of NTRUEncrypt’s private key structure to construct lower-dimensional “zero-run” lattices and “dimension-reducing” lattices from which an attacker could quickly recover an NTRUEncrypt-107 private key. Gentry [8] used a ring homomorphism from R to $\mathbb{Z}[X]/(X^{N/d} - 1)$ to “fold” L_{CS} into a more manageable lattice of dimension $2N/d$ for N having a nontrivial divisor d .

4.2 GCD Lattice Attack

Since reducing L_{CS} does indeed appear to be infeasible, the natural inclination of the cryptanalyst is to look for smaller lattices that contain the private key.

The authors of R-NSS mention one such lattice in [16]. They observe that if an attacker is able to recover the values of several $f * w$'s in R – “unreduced” modulo q – then f will likely be the shortest vector in the N -dimensional lattice formed by the rows of the several M_{f*w} 's.

Recall from section 2.2 that, for any polynomial a , there is an equivalence between the ideal (a) of a -multiples and the lattice generated by M_a . Similarly, the lattice spanned by the rows of M_{f*w_1} and M_{f*w_2} corresponds to the ideal $I = (f * w_1, f * w_2)$. Every polynomial in I is a multiple of f . Moreover, if (w_1) and (w_2) are relatively prime – i.e., there exist $a, b \in R$ with $a * w_1 + b * w_2 = 1$ – then $f \in I$, and we may say that $\text{GCD}(f * w_1, f * w_2) = (f)$. This “lattice attack” is, in fact, the standard ideal-GCD algorithm, and is among the lattice ideal operations discussed in [2].

Given how the w_i are produced in R-NSS, (f) often is indeed the GCD of two unreduced signatures, and it is even more likely to be the GCD of several unreduced signatures. Therefore, given a few unreduced signatures, we can construct an N -dimensional lattice whose shortest vector is likely f . The authors of R-NSS note that, although reducing this N -dimensional lattice is still not a trivial problem for R-NSS parameters, it is much easier than reducing the $2N$ -dimensional L_{CS} , given the exponential relationship between dimension and running time. R-NSS uses “masking” techniques to prevent recovery of $f * w$'s in R to avoid this lattice attack. (In section 5 we show that recovering $f * w$'s is nonetheless quite easy.)

4.3 Averaging Attack

The so-called “averaging attack” was first considered by Kaliski during collaboration with Hoffstein in the context of an early precursor to NSS (see patent [12]). In this work, it was observed that in the ring R , the value $f * \bar{f}$ could be obtained by an averaging attack. This attack was fatal, since in the scheme [12], f itself is a palindrome (unlike in R-NSS), thus $f * \bar{f} = f^2$. There is an efficient algorithm for taking the square root in R (see, e.g., [13]), so $f * \bar{f}$ revealed f .

In [16], the authors of R-NSS do mention this averaging attack, but also remark that knowledge of $f * \bar{f}$ does not appear to be useful. See [16], [9], [21] for a discussion of this and other ways in which the attacker may average a transcript of signatures in such a way as to get information about the private key.

Here is a description of how the averaging attack works. Suppose we can obtain a set of unreduced $f * w$'s in R . Now, consider the average

$$A_r = (1/r) \sum_{i=1}^r (f * \bar{f}) * (w_i * \bar{w}_i) \dots$$

For each i , $(w_i * \bar{w}_i)_0 = \|w_i\|^2$, which is a large positive quantity. However, for $k \neq 0$, $(w_i * \bar{w}_i)_k$ has a random distribution of positive and negative quantities that averaging essentially cancels out. Thus, as r increases, A_r essentially converges to a scalar multiple of $f * \bar{f}$. This convergence is quite fast: after a few thousand

signatures a close estimate of $f * \bar{f}$ can be computed, meaning that we obtain an estimate z such that for most coefficients, $|z_i - (f * \bar{f})_i| \leq 2$. Even if reduced signatures are used, with some corrections, there is still a convergence to $f * \bar{f}$, albeit about 10 times as slow. Clearly, the more signatures, the better the estimate. In section 6, we explain how this averaging attack may be combined with a lattice attack to recover $f * \bar{f}$ quickly and completely.

5 Lifting the Signatures

In this section, we present our first (and arguably the most important) attack with which we obtain R -multiples of the private key f . More specifically, we assume that, as passive adversaries, we are given a transcript of legitimate signatures $\{(m_1, s_1), \dots, (m_r, s_r)\}$. Using this transcript and the signer's public key, we directly compute the following elements in R_q :

$$\{f * w_1 \bmod q, \dots, f * w_r \bmod q\} \text{ and } \{g * w_1 \bmod q, \dots, g * w_r \bmod q\},$$

where the w_i are computed according to the signing process above. We will then lift these signatures to the ring R , obtaining a list of multiples of f and g :

$$\{f * w_1, \dots, f * w_r\} \text{ and } \{g * w_1, \dots, g * w_r\},$$

which are “unreduced” modulo q . This is a devastating attack against R-NSS, because undoing the q modular reduction of the signatures allows us to use the N -dimensional GCD lattice attack, described in section 4.2, rather than the $2N$ -dimensional Coppersmith-Shamir attack, reducing the key recovery time by a factor exponential in N . It also permits the other, more efficient attacks discussed later in this paper.

5.1 The Principle

From the signing procedure in the ring R_q , we get the following equations for each i :

$$f * w_i \equiv s_i \bmod q \text{ and } g * w_i \equiv t_i \bmod q, \quad (6)$$

$$f * w_i \equiv g * w_i \equiv m_i + e_i \bmod p. \quad (7)$$

If we knew e_i , we would be able to compute $f * w_i$ and $g * w_i$ modulo pq via the Chinese Remainder Theorem. We could then recover $f * w_i$ and $g * w_i$ over R without too much difficulty, since almost all of the coefficients of $f * w_i$ and $g * w_i$ will lie in the interval $(-pq/2, pq/2]$.²

The use of e_i in the signing process makes lifting the signatures much less straightforward. Since e_i will have about 20 nonzero coefficients for the suggested parameters of R-NSS, we cannot simply guess e_i from among $\binom{251}{20} 2^{20}$

² See [7] for an analysis of the coefficient distributions of convolution products. The key points here are that L2 norms $\|f\|$, $\|g\|$ and $\|w_i\|$ are small, $\|f * w_i\| \approx \|f\| \|w_i\|$ and the coefficients of $f * w_i$ have a roughly normal distribution.

possibilities. Later, we will mention how specific properties of the signing process make some possibilities much more probable than others, but even with these refinements the guessing approach remains infeasible.

Instead, we will use an iterative approach that, at each step, attempts to improve upon previous approximations. Let S_i denote our approximation of $f * w_i$, and T_i our approximation of $g * w_i$. We initialize S_i and T_i to be the polynomials in R_{pq} that satisfy

$$S_i \equiv s_i \pmod{q} \quad \text{and} \quad T_i \equiv t_i \pmod{q}, \quad (8)$$

$$S_i \equiv T_i \equiv m_i \pmod{p}. \quad (9)$$

Our approximations will have two different types of errors. First, the k th coefficients of S_i and T_i will be wrong if the k th coefficient of e_i is nonzero. Second, a coefficient of S_i (resp. T_i) may be correct modulo pq but incorrect in R if the corresponding coefficient of $f * w_i$ (resp. $g * w_i$) lies outside the interval $(-pq/2, pq/2]$. On average, about 25 (out of 251) coefficients of our initial approximations will be incorrect.

In refining our approximations, we begin with the following observation: For any (i, j) ,

$$(f * w_i) * (g * w_j) - (f * w_j) * (g * w_i) = 0. \quad (10)$$

Based on this observation, we would expect $S_i * T_j - S_j * T_i$ to tend towards 0 as our approximations improve. In fact, this is the case. We can use the norms $n_{ij} = \|S_i * T_j - S_j * T_i\|$ to decide what adjustments we should make, and to know when our approximations are finally correct. The rest of the lifting procedure is simply an elaboration of this basic idea, and one can imagine a variety of different methods through which an attacker could use these norms to create an effective lifting procedure. In our particular implementation against R-NSS, we used the norms n_{ij} , together with some R-NSS-specific heuristics, to create a very fast lifting procedure that worked almost all the time with a transcript of only four signatures.

5.2 Our Implementation of the Lifting Procedure

For each approximation pair (S_i, T_i) , we computed a “norm product” with the other approximation pairs according to the formula $P_i = \prod_{j \neq i} n_{ij}$. Preferring approximation pairs with higher norm products, we then picked a random pair (S_i, T_i) to be corrected. For each coefficient position, we temporarily added or subtracted certain multiples of q to S_i and T_i (since the approximations are already correct modulo q), and recomputed P_i . With a little bookkeeping, this step can be made extremely fast. We preserved the adjustment that reduced P_i by the greatest amount. Finally, we terminated this process when the norm product of some approximation pair reached zero, at which point we would have two correct approximation pairs.

We also used some heuristics based on specific properties of the signing process to improve the performance. We make the observation that the coefficients

of $(f * w_i, g * w_i)$ come from a probability distribution, which by the equations above, is correlated to that of e_i and the initial value of (S_i, T_i) . Specifically, the way in which e_i is computed makes it highly probable that if the k th coefficient of e_i is nonzero, then the k th coefficients of the initial approximation pair (S_i, T_i) will both be fairly small – e.g., in the range $[-100, 100]$. On the other hand, if a coefficient of the initial S_i or T_i is off by a multiple of pq , then it is highly probable that this coefficient is fairly large – e.g., outside the range $[-100, 100]$. Consequently, if corresponding coefficients of the initial S_i and T_i are both in $[-100, 100]$, we might only consider adjusting these coefficients by $\pm q$; and if either of these coefficients is outside $[-100, 100]$, we might only try adjusting by $-pq$ or pq (depending on whether the original coefficient was initially positive or negative, respectively). In our implementation, we used precisely this heuristic in the early stages of our lifting procedure. We would then “shift gears,” considering adjustments by other multiples of q so that we could catch low-probability errors as well.

We note that not every adjustment made during the lifting procedure is actually a correction. Often, this procedure will make a previously correct coefficient incorrect, and then switch it back later on. In other words, it behaves somewhat like a “random walk”. This fact, together with the heuristic nature of the overall algorithm, admittedly makes the lifting procedure difficult to analyze. For the specified parameters of R-NSS, however, it works quickly and reliably. For a transcript of four signatures, it is able to lift two signatures 90% of the time in an average of about 25 seconds (on a desktop computer). In the remaining 10%, the number of errors never converges to zero. For three signatures, it still works 70% of the time, typically finishing in about 15 seconds.

6 Obtaining $f * \bar{f}$

An important ingredient of the final algorithm is the product of f with its reversal, \bar{f} . In order to recover $f * \bar{f}$ in the context of R-NSS, we used a combination of the averaging attack mentioned in section 4.3 and a lattice attack on $f * \bar{f}$ noticed by the authors and Jonsson, Nguyen and Stern.

The lattice attack is a derivative of the Coppersmith-Shamir attack described in section 4.2. Since sending a polynomial to its reversal is an automorphism, $(f * \bar{f}) * (h * \bar{h}) \equiv (g * \bar{g}) \pmod{q}$. This means that the vector $(f * \bar{f} \| g * \bar{g})$ is contained in the lattice L_{norm} generated by

$$B_{norm} = \begin{bmatrix} I_{(N)} & M_{h*\bar{h}} \\ 0 & qI_{(N)} \end{bmatrix}. \quad (11)$$

This lattice has dimension $2N$, but it has an $(N + 1)$ -dimensional sublattice of palindromes, which contains $(f * \bar{f} \| g * \bar{g})$. Conceivably, recovering $(f * \bar{f} \| g * \bar{g})$ from this sublattice could give us useful information about f and g . However,

this attack fails for typical NTRU or NSS parameters, since $(f * \bar{f} \| g * \bar{g})$ is normally not the shortest vector.³

For R-NSS, we combine ideas from the above attack with the GCD attack in 4.2 and the averaging technique in 4.3. First, we use our unreduced signatures to form the ideal $(f * \bar{f})$ from a few unreduced signature products $s * \bar{s} = f * \bar{f} * w_i * \bar{w}_i$, exactly as described in Section 4.2. Then, we take the subring of $(f * \bar{f})$ consisting of palindromes, which forms a lattice of dimension $(N + 1)/2$. In fact, this lattice is generated by $f * \bar{f}$, and $(N - 1)/2$ vectors $(X^k + X^{N-k}) * f * \bar{f}$. For the same reason as above, $f * \bar{f}$ might not be the shortest vector in the lattice. However, we may use the averaging attack to obtain a good estimate t of $f * \bar{f}$, modify the lattice to include t , and then use lattice reduction to obtain the (shortest) vector $t - f * \bar{f}$.⁴ In practice, this attack is amazingly effective for two reasons: the lattice problem is only $(N + 1)/2$ dimensional, and $\|t - f * \bar{f}\|$ will be much less than $\|f * \bar{f}\|$ for even a *very* poor estimate of t . We found that we needed only 10 signatures to obtain a sufficiently accurate estimate t of $f * \bar{f}$ (even though only a handful of coefficients in t were actually correct). With these 10 signatures, we consistently recovered $f * \bar{f}$ in less than 10 seconds.

7 Orthogonal Congruence Attack

In this section, we describe a polynomial-time algorithm for recovering the private key f from $f * \bar{f}$ and one other multiple of f , such as $f * \bar{w}$, when w is relatively prime to \bar{f} . In other words, this algorithm requires $f * \bar{f}$ and a basis B_f of the ideal (f) .⁵ This algorithm is quite surprising, and uses novel ideas combining orthogonal lattices with number theoretic congruence arising from the cyclotomic field $\mathbb{Q}(\zeta_N)$.

The complete algorithm is rather complex, but here is a brief (and not entirely accurate) sketch: We begin by choosing a large prime number $P \equiv 1 \pmod{N}$. (For now, we defer discussing how large P must be.) Then, using $f * \bar{f}$ and our basis for (f) , we use a series of lattice reductions to obtain $f^{P-1} * a$ for some polynomial a , and a guarantee that $\|a\| < P/2$. Using the congruence $f^{P-1} \equiv 1 \pmod{P}$, we will be able to compute $a \pmod{P}$ and hence a exactly, from which we will be able to compute f^{P-1} exactly. We will then use this power of f to recover f .

³ By the ‘‘Gaussian heuristic,’’ the expected length of the shortest vector in a lattice of determinant d and dimension n is $d^{1/n} \sqrt{n/(2\pi e)}$. For L_{norm} , this length is $\sqrt{qN}/(\pi e)$. On the other hand, since $\|f\| > \sqrt{N}$ in NSS and $\|f * \bar{f}\| \geq \|f\|^2$ (the latter inequality following from $(f * \bar{f})_0 = \|f\|^2$), the norm of $f * \bar{f}$ is greater than N and hence greater than the Gaussian heuristic, since N is typically chosen to be greater than q .

⁴ One could also use Babai’s algorithm to solve the closest vector problem (CVP).

⁵ Yet another characterization of the algorithm is that it recovers f from B_f and the relative norm of f over the index 2 subfield of $\mathbb{Q}(\zeta_N)$.

We describe this algorithm and the theory behind it in more detail below. Our first task will be to find a tool that ensures that when LLL gives us $f^{P-1} * a$, there is a definite bound on $\|a\|$.

7.1 Orthogonal Lattices

Certain lattices possess a basis of N equal length, mutually perpendicular basis vectors. We denote such lattices *orthogonal lattices*. Two lattices are called *homothetic* if up to a constant stretching factor, λ , there is a distance preserving map from one lattice to the other. That is, all orthogonal lattices are homothetic to the trivial lattice \mathbb{Z}^N . Similarly, we define f to be an *orthogonal polynomial* if the circulant matrix M_f is the orthogonal basis of an orthogonal lattice. We are interested in orthogonal lattices because they possess a multiplicative norm property.

We note that for randomly chosen polynomials a and f , the norm is quasi-multiplicative, $\|f * a\| \approx \|f\| \cdot \|a\|$. However, if one of the polynomial factors, say f , is *orthogonal*, then equality will hold

$$\|f * a\| = \|f\| \cdot \|a\| . \quad (12)$$

For general polynomials f , applying LLL to (f) is guaranteed to find a multiple of f , say $f * a$, such that the norm $\|f * a\|$ is less than a specific factor times the norm of the shortest vector in the lattice. In the case where f is orthogonal, we can additionally bound $\|a\|$ by this factor, since $\|f * a\| = \|f\| \cdot \|a\|$. In the case of LLL, this means that we can be certain that $\|a\| < 2^{(N-1)/2}$.

7.2 Using $f * \bar{f}$ to Construct an Implicit Orthogonal Lattice

What do we do when f is not an orthogonal polynomial, but our objective is to find $f * a$ with small $\|a\|$ (given only $f * \bar{f}$ and lattice basis B_f of (f))? Of course, we may apply LLL to B_f and just hope that the output vector $f * a$ has short a , but this may not work even if f is only slightly nonorthogonal⁶. This section describes how we can accomplish this task by using knowledge of $f * \bar{f}$ to *implicitly* define an orthogonal lattice.

Since B_f and M_f are both bases of (f) , they are related by $B_f = U \cdot M_f$ for some unimodular matrix U . Notice that each row of B_f is of the form $f * u_i$ where u_i is the i th row of U . This means that the objective of finding $f * a$ with bounded $\|a\|$ is equivalent to bounding the norms of the rows of U . So, in some sense, we would like to apply lattice reduction to U . How can we reduce the rows of U when we only know $B_f = U \cdot M_f$, and not U itself?

Supposing that f is not a zero divisor in R , we can also divide by $f * \bar{f}$. Allowing denominators in our notation, we let $D = M_{(1/(f*\bar{f}))}$ and compute

$$B_f \cdot D \cdot B_f^T = U \cdot U^T , \quad (13)$$

⁶ The notion of nonorthogonality is made precise with concept called *orthogonality defect*.

which is the Gram matrix of our unknown unimodular matrix U . Although we do not know U explicitly, $U \cdot U^T$ has all the information that LLL needs to know about U in order to reduce it – namely, the mutual dot products $u_i \cdot u_j$ of each pair of row vectors. We can therefore apply a Gram matrix version of LLL to $U \cdot U^T$, which outputs the unimodular transformation matrix V , and the Gram matrix of the reduced lattice: $(V \cdot U) \cdot (V \cdot U)^T$. By the LLL bound, the norms of the rows of (the unknown) basis $V \cdot U$ will be bounded by $2^{(N-1)/2}$. Now, we can compute a new basis of (f) – namely, $(V \cdot U) \cdot M_f = V \cdot B_f$ – and be certain that each row of this basis equals $f * a_i$ for $\|a_i\| < 2^{(N-1)/2}$. Effectively, we have reduced the orthogonal lattice defined by U , without even knowing an explicit basis for it.

7.3 Galois Congruence

In addition to the orthogonal lattices technique, we use some interesting congruences on the ring R . The first congruence states that for any prime P such that $P \equiv 1 \pmod{N}$,

$$f^P = f \pmod{P} . \quad (14)$$

This implies that for any f which is not a zero divisor⁷ in $R_P = \mathbb{Z}_P[X]/(X^N - 1)$ that

$$f^{P-1} = 1 \pmod{P} . \quad (15)$$

We may generalize these equations to arbitrary primes P by using a *Galois Conjugation* function (written as a superscript) $\sigma(r) : R \rightarrow R$, defined by

$$f^{\sigma(r)}(x) = f(x^r) . \quad (16)$$

For any r not divisible by N , $\sigma(r)$ defines an automorphism on R . There are $N-1$ such automorphisms, since two values of r which differ by a factor of N define the same automorphism. We call $\sigma(r)$ the r th Galois conjugation mapping, and have the congruence

$$f^P = f^{\sigma(P)} \pmod{P} . \quad (17)$$

For elementary proofs of equations 15 and 17 and their relationship with the the Galois theory of $Q(\zeta)$, we refer the reader to the Appendix E; for the relationship with the the Galois theory of $Q(\zeta)$, see Appendix D.

As mentioned above, our motivation for considering such congruences is that given a multiple of f^{P-1} , say $f^{P-1} * a$, we may use the congruence $f^{P-1} = 1 \pmod{P}$ to conclude that

$$f^{P-1} * a = a \pmod{P} . \quad (18)$$

Now, if a is so small that all of its coefficients lie in the interval $(-P/2, P/2]$, then the representatives for $f^{P-1} * a \pmod{p}$ in this interval reveal a exactly. Assuming that a is not a zero divisor in R , dividing the product by a then

⁷ See Appendix E for a discussion of the zero divisor issue in R_P .

yields the exact value of f^{P-1} . With this observation, we are in a position to use orthogonal lattice theory with lattice reduction to obtain small multiples of powers of f and thus exact powers of f .

However, there is a technical difficulty arising from the fact that LLL only guarantees that $\|a\| < 2^{(N-1)/2}$. To ensure that a has coefficients in $(-P/2, P/2]$, P has to be quite large – about the same order of magnitude as $2^{(N-1)/2}$. This means that f^{P-1} has bit-length exponential in N , which makes it impossible for us to even store the value of this polynomial. Initially, this might appear to be a fatal problem. It also indicates that the “brief sketch” given at the beginning of Section 7 could not have been entirely accurate.

Fortunately, we will not need to work with f^{P-1} directly; it will be sufficient for our purposes to be able to compute f^{P-1} modulo some primes. As discussed further in section 7.4, we can make such computations using a process similar to repeated squaring. However, to recover f , we will need some power of f that we can work with directly (unreduced). To solve this problem, we may choose a second prime $P' \equiv 1 \pmod{N}$ with $\text{GCD}(P-1, P'-1) = 2N$ for which we can compute $f^{P'-1}$ modulo some primes. Then, using the Euclidean Algorithm, we may compute f^{2N} modulo these primes, and ultimately f^{2N} exactly via the Chinese Remainder Theorem. We may work with f^{2N} directly, since its bit-length is merely polynomial in N . A more elegant solution in Appendix F computes f^{2N} directly from f^{P-1} if $\text{GCD}(P-1, 2N-1) = 1$. In Section 7.5, we describe how to recover f from f^{2N} .

7.4 Ideal Power Algorithms

In practice, it might be the case that smaller P would be sufficient, allowing us to work with the ideal (f^{P-1}) directly, but some tricks are needed to handle the very large primes P that the LLL bound imposes on us. Suppose we had the chains of polynomials $\{v_0^2 * \overline{v_1}, \dots, v_{r-1}^2 * \overline{v_r}\}$ and $\{v_0 * \overline{v_0}, \dots, v_{r-1} * \overline{v_{r-1}}\}$. Then, we could make the following series of computations:

$$v_0^4 * \overline{v_2} = (v_0^2 * \overline{v_1})^2 * (v_1 * \overline{v_1})^{-2} * (v_1^2 * \overline{v_2}) \pmod{P}, \quad (19)$$

$$v_0^8 * \overline{v_3} = (v_0^4 * \overline{v_2})^2 * (v_2 * \overline{v_2})^{-2} * (v_2^2 * \overline{v_3}) \pmod{P}, \quad (20)$$

and so on, ending with the equation:

$$v_0^{2^r} * \overline{v_r} = (v_0^{2^{r-1}} * \overline{v_{r-1}})^2 * (v_{r-1} * \overline{v_{r-1}})^{-2} * (v_{r-1}^2 * \overline{v_r}) \pmod{P}. \quad (21)$$

In other words, we could compute $v_0^{2^r} * \overline{v_r} \pmod{P}$ efficiently even though the exponent 2^r may be quite large. If we could use this approach to get $v_0^{P-1} * \overline{v_r} \pmod{P}$ where $\|\overline{v_r}\| < P/2$, then we could recover $\overline{v_r}$ exactly, and then we could use the same chains of polynomials to compute v_0^{P-1} modulo other primes.

In this section, we describe how to get such chains of polynomials so that we can get modular data about f^{P-1} from $f * \overline{f}$ and a basis B_f of (f) . The main tool that we use is the multiplication of ideals (see [2]), and we adapt this technique to use the orthogonal lattice theory above. The algorithm described

in the paragraph above is essentially a repeated squaring algorithm, so we first review how to multiply ideals, and in particular, obtain the ideal (f^2) from the ideal (f) .

Remark 2. Ideal multiplication in R : If $A = (f)$ and $B = (g)$ are principal ideals generated as \mathbb{Z} -modules by $(f * a_1, \dots, f * a_n)$ and $(g * b_1, \dots, g * b_n)$, then the ideal product AB is generated as a \mathbb{Z} -module by the n^2 elements of the set $\{a_i * b_j * f * g\}$, which defines $(f * g)$.

Note that we describe the ideals as modules – i.e., as the \mathbb{Z} -span of a set of polynomials (rather than by giving generators over R). This is because our algorithms represent ideals as lattices – i.e., as lists of polynomials. We use ideal multiplication as follows: Given the ideal (f) in terms of the basis $B_f = U \cdot M_f$, we can generate (f^2) from the rows $b_i * b_j = u_i * u_j * f^2$. Since we know $f^2 * \bar{f}^2$, we can use the orthogonal lattice method described in section 7.2 to obtain a reduced basis $B_{f^2} = U' M_{f^2}$ where the rows of U' have norm less than $2^{(N-1)/2}$.⁸ Next, we pick a row of B_{f^2} and name it $f^2 * \bar{v}_1$. We can directly compute $v_1 * \bar{v}_1$ and a basis for $v_1: U' \cdot M_{v_1}$. At this point we can compute a basis for (v_1^2) to begin another iteration of this process. Thus, we have the chains of polynomials previously introduced, and, with a modification to multiply some of the ideals by (f) , its natural generalization. (Note: v_0 represents f in the following theorem.)

Theorem 1. *Polynomial Chains*

Suppose v_0 is not a zero divisor in $R = \mathbb{Z}[X]/(X^N - 1)$. Let $k = \sum k_i 2^i$ with $k_i \in \{0, 1\}$ be an integer with $r = \lfloor \log_2(k) \rfloor$. Let P be a prime such that v_0 is not a zero divisor in R_P . Then, given the input $v_0 * \bar{v}_0$ and basis B_0 of (v_0) , we may compute, in time polynomial in r , N and the bit-length of the input, the chains:

$$\{v_0^{k_{r-1}} * v_0^2 * \bar{v}_1, \dots, v_0^{k_0} * v_{r-1}^2 * \bar{v}_r\} \quad \text{and} \\ \{v_0 * \bar{v}_0, \dots, v_{r-1} * \bar{v}_{r-1}\},$$

where for $\forall i > 0$, no v_i is a zero divisor in R_P , and $\|v_i\| < 2^{(N-1)/2}$. Using these chains, we may compute $v_0^k * \bar{v}_r \pmod{P}$ in polynomial time. If $k = P - 1 \geq 2^{(N+1)/2}$ with $P \equiv 1 \pmod{N}$, we may compute \bar{v}_r exactly, and thereafter use the above chains to compute $v_0^{P-1} \pmod{Q}$ in polynomial time for any prime Q such that no v_i is zero divisor in R_Q .

Notice that the above chains are identical to those given in the toy example at the beginning of this section, except that k is permitted not to be a power of 2. Zero-divisor issues arise because, at times, certain polynomials are divided out, as in the toy example, but all zero divisor issues are easily dealt with by using determinant-based technique noted in Appendices E and F. For $i > 0$, we have $\|v_i\| < 2^{(N-1)/2}$ by the orthogonal lattice method described in section 7.2. This

⁸ In the worst case, this involves reduction of an N -dimensional lattice defined by N^2 generators. This reduction would only have polynomial time-complexity, but normally we will be able to do much better, since we will usually be easy to find far fewer (on the order of N) polynomials that span the ideal.

norm bound and the polynomial running time of LLL underlie the polynomial complexity of the algorithm. In the Appendix F, we write this algorithm in terms of pseudocode in an effort to clarify its details as well as its polynomial-time complexity. Now, we will use the algorithm embodied in Theorem 1 to obtain f^{2N} .

Theorem 2. *Computing f^{2N}*

Assume f is not a zero divisor in $R = \mathbb{Z}[X]/(X^N - 1)$. Then, given $f * \bar{f}$ and basis B_f of (f) , we may compute f^{2N} in time polynomial in N and the bit-length of f .

Proof. We choose primes P and P' , each greater than $2^{(N+1)/2}$, with $\text{GCD}(P - 1, P' - 1) = 2N$ and f is a zero divisor in neither R_P nor $R_{P'}$ (using Dirichlet's theorem on primes in arithmetic progression and the fact that f may be a zero divisor in R_Q for only a finite number of primes Q). By Theorem 1, we may compute two polynomial chains that will allow us to compute $f^{P-1} \pmod{r_i}$ and $f^{P'-1} \pmod{r_i}$ in polynomial time for any prime r_i such that no v_i (in either chain) is zero divisor in R_{r_i} . (We simply avoid the finite number of problematic primes.) Applying the Euclidean algorithm, we compute f^{2N} modulo each r_i , and ultimately f^{2N} exactly using the Chinese Remainder Theorem.

7.5 Computing f from f^{2N}

Our final task is to compute the private key f from f^{2N} . We use the following theorem.

Theorem 3. *Galois Polynomial*

Given f^{2N} and $b \in \mathbb{Z}_N^*$, we may compute $z = f^{\sigma(-b)} f^b$ in time polynomial in b , N and bit-length of f .

Proof. Let $P_2 > 2\|f^{\sigma(-b)} f^b\|$ be a prime number such that $P_2 = 2c_2N - b$ for some integer c_2 . Then, $(f^{2N})^{c_2} \equiv f^{P_2} * f^b \equiv f^{\sigma(-b)} f^b \pmod{P_2}$. Since $P_2 > 2\|f^{\sigma(-b)} f^b\|$, we recover $z = f^{\sigma(-b)} f^b$ exactly.

Now, in terms of recovering f , we first note that f^{2N} uniquely defines f only up to sign and rotation – i.e., up to multiplication by $\pm X^k$, the $2N$ th roots of unity in R . The basic idea of our approach is that, given f^{2N} , fixing $f(\zeta)$ for one (complex) N th root of unity ζ completely determines $f(\zeta^d)$ for all exponents d . Then, we may use the $N - 1$ values of $f(\zeta^d)$, together with $f(1)$ (which we will know up to sign), to solve for f using Gaussian elimination. If we set $-b$ to be a primitive root modulo N , the polynomial z given in Theorem 3 will help us iteratively derive the $f(\zeta^d)$ from $f(\zeta)$ as follows:

$$f(\zeta^{(-b)^{i+1}}) = z(\zeta^{(-b)^i}) / f^b(\zeta^{(-b)^i}) . \tag{22}$$

Repeated exponentiation will not result in a loss of precision, since the value may be corrected at each stage. Since $-b$ is a primitive root modulo N , these evaluations give us $N - 1$ linearly independent equations in the coefficients of f , which together with $f(1)$, allow us to recover the private key f completely, up to sign and rotation.

8 Summary and Generalizations of R-NSS Cryptanalysis

For the reader’s convenience, we briefly review the main points of the attack. The first two stages of the attack are fast in practice, but they are both heuristic and have no proven time bounds. The lifting procedure lifts a transcript of signatures from R_q to R , obtaining unreduced f -multiples in R . The second stage uses an averaging attack to approximate $f * \bar{f}$, and then solves the closest vector problem (CVP) to recover $f * \bar{f}$ exactly. The algorithm of the final stage, which we have not fully implemented, uses output from the previous two stages to recover the private key in polynomial time. By combining lattice-based methods and number-theoretic congruences, the algorithm of the final stage can be used to:

1. Recover f from $f * \bar{f}$ and a basis B_f of (f) ;
2. Recover f from only B_f when f is an orthogonal polynomial; and
3. Recover f/\bar{f} from B_f whether f is an orthogonal polynomial or not.

We anticipate that this algorithm could be generalized to recover f given a basis of (f) and the relative norm of f over an index 2 subfield where the degree-2 extension is complex conjugation. In Section 9, we discuss another possible generalization of this algorithm that may be an interesting area of research.

9 NTRUSign

NTRUSign was proposed at the rump session of Asiacrypt 2001 as a replacement of R-NSS [11], and, as requested by the Program Committee, we provide some preliminary security analysis. The scheme is more natural than previous NTRU-based signature schemes, particularly in terms of its sole verification criterion: the signer (or forger) must solve an “approximate CVP problem” in the NTRU lattice – i.e., produce a lattice point that is sufficiently close to a message digest point, à la Goldreich, Goldwasser and Halevi [10]. Similar to the GGH cryptosystem, the signer has private knowledge of a “good” basis of the NTRU lattice having short basis vectors, and publishes the usual “bad” NTRU lattice basis:

$$B_{priv} = \begin{bmatrix} M_f & M_g \\ M_F & M_G \end{bmatrix} \quad B_{pub} = \begin{bmatrix} I_{(N)} & M_h \\ 0 & qI_{(N)} \end{bmatrix} .$$

Unlike GGH, these bases may be succinctly represented with only four polynomials: (f, g, F, G) for the private basis, and h for the public basis. (Recall that M_f denotes the circulant matrix corresponding to the polynomial f ; see Section 2.) In terms of key generation, the signer first generates short polynomials f and g and computes the public key as $h = f^{-1} * g \pmod{q}$, as in NTRUEncrypt or R-NSS. Due to lack of space, we refer the reader to [11] for details on how the signer generates his second pair of short polynomials (F, G) , but we note the following properties: 1) $f * G - g * F = q$ and 2) $\|F\|$ and $\|G\|$ are 2 to 3 times greater than $\|f\|$ and $\|g\|$.

To sign, the message is hashed to create a random message digest vector (m_1, m_2) with $m_1, m_2 \in R_q$. The signer then computes:

$$G * m_1 - F * m_2 = A + q * C , \quad (23a)$$

$$-g * m_1 + f * m_2 = a + q * c , \quad (23b)$$

where A and a have coefficients in $(-q/2, q/2]$, and sends his signature:

$$s \equiv f * C + F * c \pmod{q} . \quad (24)$$

The verifier computes $t \equiv s * h \pmod{q}$ and checks that (s, t) is “close enough” to (m_1, m_2) – specifically,

$$\|s - m_1\|^2 + \|t - m_2\|^2 \leq \text{Normbound} . \quad (25)$$

We can see why verification works when we write, say, s in terms of Equations 23a and 23b:

$$s \equiv m_1 - (A * f + a * F)/q \pmod{q} , \quad (26)$$

where $\|A * f + a * F\|$ will be reasonably short since f and F are short.

In the absence of a transcript, the forgery problem is provably as difficult as the approximate CVP problem in the NTRU lattice. However, it is clear that NTRUSign signatures leak some information about the private key. The mapping involution $\mathbb{Z}[X]/(q, X^N - 1)$ sending $m \mapsto s$ is not a permutation, and the associated identification protocol is not zero knowledge (even statistically or computationally). Below, we describe concrete transcript attacks using ideas from our cryptanalysis of R-NSS. We have had fruitful discussions with NTRU regarding these attacks, and they have begun running preliminary tests to determine their efficacy. Based on these tests, some of these attacks appear to require a very long transcript that may make them infeasible in practice. This is a subject of further research. In any case, these attacks show that NTRUSign cannot have any formal security property, since it is not secure against passive adversaries.

9.1 Second Order Attack

Using Equation 26, we may obtain polynomials of the form $(A * f + a * F)$ (similarly, $(A * g + a * G)$). Consider the following average:

$$\text{Avg}_{ff}(r) = (1/r) \sum_{i=1}^r (a_i * F + A_i * f) * \overline{(a_i * F + A_i * f)} \quad (27)$$

$$= (1/r) \sum_{i=1}^r (a_i * \overline{a_i}) * (F * \overline{F}) + (A_i * \overline{A_i}) * (f * \overline{f}) + \text{other terms} . \quad (28)$$

The “other terms” will converge to 0, since A and a are uniformly distributed at random modulo q and, though dependent, have small statistical correlation. (See

Remark 3 below.) The explicit portion of the average will converge essentially to a scalar multiple of $f * \bar{f} + F * \bar{F}$, for the same reasons as discussed in Section 4.3. Thus,

$$\lim Avg_{ff}(r) = \gamma(f * \bar{f} + F * \bar{F}) . \quad (29)$$

Because the signatures in a transcript are random variables, the limit converges as $1/\sqrt{r}$ where r is the length of a transcript. We may use this averaging to obtain a sufficiently close approximation of $f * \bar{f} + F * \bar{F}$ to obtain the exact value by solving the CVP in the $(N + 1)$ -dimensional lattice given in Equation 11 using lattice reduction. Thus, we recover a polynomial that is quadratic in the private key, and we can obtain $f * \bar{g} + F * \bar{G}$ and $g * \bar{g} + G * \bar{G}$ in a similar fashion.

Remark 3. One may artificially construct situations where $(1/r) \sum_{i=1}^r a_i * \bar{A}_i$ does not converge to 0. For example, if we let $f = \bar{F}$, then $A_i * f + a_i * F \equiv 0 \pmod{q}$ basically implies $A_i = -a_i$ and hence $a_i * \bar{A}_i = -a_i * \bar{a}_i$, the average of which does not converge to 0. Conceivably, NTRUSign could be modified so as to constrain $f^{-1} * F \pmod{q}$, but this would likely allow alternative attacks.

It is worth noting that these second order polynomials give us the Gram matrix $B_{priv}^T \cdot B_{priv}$:

$$B_{priv}^T \cdot B_{priv} = \begin{bmatrix} M_{\bar{f}} & M_{\bar{F}} \\ M_{\bar{g}} & M_{\bar{G}} \end{bmatrix} \begin{bmatrix} M_f & M_g \\ M_F & M_G \end{bmatrix} = \begin{bmatrix} M_{f*\bar{f}+F*\bar{F}} & M_{g*\bar{f}+G*\bar{F}} \\ M_{f*\bar{g}+F*\bar{G}} & M_{g*\bar{g}+G*\bar{G}} \end{bmatrix} .$$

This Gram matrix gives us the “shape” of the parallelepiped defined by B_{priv}^T , but the “orientation” of this parallelepiped is unclear. An interesting (and open) question is: Can an attacker recover B_{priv} from $B_{priv}^T \cdot B_{priv}$ and $B_{pub} = U \cdot B_{priv}$, where U is a unimodular matrix? We answered a similar question in the affirmative in Section 7; we showed that an attacker, in polynomial time, can recover M_f from $M_f^T \cdot M_f$ and $U \cdot M_f$, where U is a unimodular matrix. We have not found a way to extend the orthogonal congruence attack to solve the NTRUSign Gram matrix problem, however, where the bi-circulant (rather than purely circulant) nature of the matrices in question (such as B_{priv}) destroys the commutativity that our orthogonal congruence attack appears to require, but this does not imply that the NTRUSign Gram matrix problem is necessarily hard. We note that it would more than suffice to find an algorithm that factors $U \cdot U^T$ for unimodular U . (This factorization is unique up to a signed permutation matrix.) Further research in this area would be interesting, if only because it is relevant to NTRUSign’s security.

9.2 Second Order Subtranscript Attack

It is clear that the second order polynomials recovered above contain information not contained in the public key, but using this information to create an effective attack is not so straightforward. Our approach has been to use the second order polynomials to recover, say, $f * \bar{f}$, so that we may then apply the orthogonal

congruence attack to recover f . One way to get $f * \bar{f}$ is to use the following subtranscript attack.

First, we notice that since $\|F\| > \|f\|$, the norm $\|A * f + a * F\|$ is dictated more by $\|a\|$ than by $\|A\|$. More relevantly, for our purposes, an $A * f + a * F$ that is longer than normal will usually have $\|a\| > \|A\|$. This suggests a subtranscript attack, including in Equation 27 only those polynomials $A_i * f + a_i * F$ for which $\|A_i * f + a_i * F\|$ is greater than some bound.⁹ Then, we have:

$$\text{Subtranscript: } \lim_{r \rightarrow \infty} \text{Avg}_{ff}(r) = \gamma_1(f * \bar{f}) + \gamma_2(F * \bar{F}), \quad (30)$$

for $\gamma_1 < \gamma_2$. Since this linear combination of $f * \bar{f}$ and $F * \bar{F}$ is distinct from that in Equation 29, we may compute $f * \bar{f}$ and $F * \bar{F}$. The convergence of this subtranscript averaging will be affected by the proportional size of the subtranscript, but more importantly, by the fact that γ_1 may be only a few percentage points greater than γ_2 (in our preliminary experiments using the longest 50% of the $A_i * f + a_i * F$'s). Further experiments are necessary to determine the effectiveness of this attack. Another consideration is that in [11], a possible modification of NTRUSign was proposed in which one chooses the transpose of B_{priv} to be the private key. The basis vectors are then (f, F) and (g, G) with $\|f\| \approx \|g\|$ and $\|F\| \approx \|G\|$. Choosing the private basis in this way appears to defeat this subtranscript attack.

9.3 Fourth Order Attack

An alternative way to get $f * \bar{f}$ is to use the following fourth order attack. Viewing the average in Equation 27, one may consider the corresponding variance and conclude, under the assumption of the statistical independence of a and A , that:

$$\lim_{r \rightarrow \infty} (1/r) \sum_{i=1}^r (s * s_{rev})^2 - (1/r) \beta \left(\lim_{r \rightarrow \infty} \sum_{i=1}^r (s * s_{rev}) \right)^2 = \gamma f * \bar{f} * F * \bar{F}. \quad (31)$$

The adjustment value of β depends on the scheme parameters n and q , and so the above value may not be exactly the variance. The factor γ also depends on the scheme parameters, and is a constant that slows convergence by a factor $\frac{1}{\gamma^2}$. This limit does converge more slowly than the second order averaging, but, as above, we may use a close approximation in conjunction with the lattice of Equation 11 to obtain the exact value. Preliminary tests show that it may not be practical to obtain an error lower than the Gaussian estimate with a reasonable number of signatures. Assuming we do obtain the value $f * \bar{f} * F * \bar{F}$, we may use it in combination with $(f * \bar{f} + F * \bar{F})^2$ to obtain $(f * \bar{f} - F * \bar{F})^2$, and then $f * \bar{f} - F * \bar{F}$ (using, perhaps, the algorithm given in Section 7.5). Then, $f * \bar{f} + F * \bar{F}$ and $f * \bar{f} - F * \bar{F}$ give us $f * \bar{f}$ and $F * \bar{F}$.

⁹ This selection criterion might be refined, as by also considering the norm of $A * \bar{F} - a * \bar{f}$, which may be computed using the above second order polynomials.

9.4 Preliminary Conclusion

The approach of these initial attacks was to reduce the breaking problem to the orthogonal congruence attack using the results of various averagings. We showed how this could be done, but the practical feasibility of these attacks has yet to be determined. In our experiments, we have found that the second order attack is feasible; for example, by averaging 20000 signatures, an attacker may obtain an approximation whose squared error is about 88, about 1/20 of the squared Gaussian heuristic of the $(N + 1)$ -dimensional CVP lattice that would be used to correct this approximation. Although we have not tested this CVP lattice, we believe its reduction would be feasible. Alternatively, more signatures could first be used to obtain a better approximation. We have also shown how the security of NTRUSign rests on the hardness of several new hard problems. These attacks will continue to be analyzed by the authors, NTRU corporation, and the cryptographic community.

10 Acknowledgments

The authors would like to thank Burt Kaliski, Alice Silverberg and Yiqun Lisa Yin for helpful discussions, Jakob Jonsson, Phong Nguyen, and Jacques Stern for discussions and collaboration on the precursor of this article, and Jeffrey Hoffstein, Nick Howgrave-Graham, Jill Pipher, Joseph Silverman and William Whyte, who have given us valuable feedback on our cryptanalysis, particularly on our preliminary cryptanalysis of NTRUSign.

References

1. M. Ajtai, *The shortest vector problem in L_2 is NP-hard for randomized reductions*, in Proc. 30th ACM Symposium on Theory of Computing, 1998, 10–19.
2. H. Cohen, *A Course in Computational Algebraic Number Theory*, Graduate Texts in Mathematics, 138. Springer, 1993.
3. H. Cohen, *Advanced Topics in Computational Number Theory*, Graduate Texts in Mathematics 138, 1993.
4. Consortium for Efficient Embedded Security. Efficient Embedded Security Standard (EESS) # 1: Draft 1.0. Previously on <http://www.ceesstandards.org>.
5. Consortium for Efficient Embedded Security. Efficient Embedded Security Standard (EESS) # 1: Draft 2.0. Previously on <http://www.ceesstandards.org>.
6. Consortium for Efficient Embedded Security. Efficient Embedded Security Standard (EESS) # 1: Draft 3.0. Available from <http://www.ceesstandards.org>.
7. D. Coppersmith and A. Shamir, *Lattice Attacks on NTRU*, in Proc. of Eurocrypt '97, LNCS 1233, pages 52–61. Springer-Verlag, 1997.
8. C. Gentry, *Key Recovery and Message Attacks on NTRU-Composite*, in Proc. of Eurocrypt '01, LNCS 2045, pages 182–194. Springer-Verlag, 2001.
9. C. Gentry, J. Jonsson, J. Stern, M. Szydlo, *Cryptanalysis of the NTRU signature scheme*, in Proc. of Asiacrypt '01, LNCS 2248, pages 1–20. Springer-Verlag, 2001.

10. O. Goldreich, S. Goldwasser, S. Halevi, *Public-key Cryptography from Lattice Reduction Problems*, in Proc. of Crypto '97, LNCS 1294, pages 112–131. Springer-Verlag, 1997.
11. J. Hoffstein, N. Howgrave-Graham, J. Pipher, J.H. Silverman, W. Whyte, *NTRUSign: Digital Signatures Using the NTRU Lattice*, December, 2001. Available from <http://www.ntru.com>.
12. J. Hoffstein, B.S. Kaliski, D. Lieman, M.J.B. Robshaw, Y.L. Yin, *Secure user identification based on constrained polynomials*, US Patent 6,076,163, June 13, 2000.
13. J. Hoffstein, D. Lieman, J.H. Silverman, *Polynomial Rings and Efficient Public Key Authentication*, in Proc. International Workshop on Cryptographic Techniques and E-Commerce (CrypTEC '99), Hong Kong, (M. Blum and C.H. Lee, eds.), City University of Hong Kong Press.
14. J. Hoffstein, J. Pipher, J.H. Silverman, *Enhanced Encoding and Verification Methods for the NTRU Signature Scheme*, NTRU Technical Note #017, May 2001. Available from <http://www.ntru.com>.
15. J. Hoffstein, J. Pipher, J.H. Silverman. *Enhanced encoding and verification methods for the NTRU signature scheme (ver. 2)*, May 30, 2001. Available from <http://www.ntru.com>.
16. J. Hoffstein, J. Pipher, J.H. Silverman, *NSS: The NTRU Signature Scheme*, preprint, November 2000. Available from <http://www.ntru.com>.
17. J. Hoffstein, J. Pipher, J.H. Silverman, *NSS: The NTRU Signature Scheme*, in Proc. of Eurocrypt '01, LNCS 2045, pages 211–228. Springer-Verlag, 2001.
18. J. Hoffstein, J. Pipher, J.H. Silverman, *NSS: The NTRU Signature Scheme: Theory and Practice*, preprint, 2001. Available from <http://www.ntru.com>.
19. A.K. Lenstra, H.W. Lenstra Jr., L. Lovász, *Factoring Polynomials with Rational Coefficients*, Mathematische Ann. 261 (1982), 513–534.
20. A. May, *Cryptanalysis of NTRU-107*, preprint, 1999. Available from <http://www.informatik.uni-frankfurt.de/~alex/crypto.html>.
21. I. Mironov, *A Note on Cryptanalysis of the Preliminary Version of the NTRU Signature Scheme*, IACR preprint server, <http://eprint.iacr.org/2001/005>.
22. P. Nguyen and J. Stern, *Lattice Reduction in Cryptology: An Update*, in Proc. of Algorithm Number Theory (ANTS IV), LNCS 1838, pages 85–112. Springer-Verlag, 2000.
23. J.H. Silverman, *Estimated Breaking Times for NTRU Lattices*, NTRU Technical Note #012, March 1999. Available from <http://www.ntru.com>.
24. J.H. Silverman, *Invertibility in Truncated Polynomial Rings.*, NTRU Technical Note #009, October 1998. Available from <http://www.ntru.com>.
25. L. Washington, *Introduction to Cyclotomic Fields*, Graduate Texts in Mathematics 83, 1982.

A Original NSS and R-NSS: The Details

This section is included for completeness, so that this paper has a complete specification of R-NSS, and so the reader can compare the original NSS to R-NSS, and see why R-NSS is not susceptible to previous attacks.

A.1 Original NSS

Like R-NSS, NSS uses the parameters N , q , p , d_f , d_g and d_m , and polynomials f, g and h , which play roughly the same role as in section 3.

Key Generation: The polynomials f and g are defined as $f = f_0 + pf_1$ and $g = g_0 + pg_1$ where f_0 and g_0 are publicly known small polynomials (typically $f_0 = 1$ and $g_0 = 1 - 2X$), and h , the public key, is computed so that $f * h = g$, as in R-NSS.

Signing: A message representative m is computed and two masking polynomials are selected at random, w_2 from $\mathcal{L}(dw_2, dw_2)$ and w_1 having at most 25 nonzero coefficients. A temporary polynomial $w = m + w_1 + pw_2$ is computed, and w_2 is altered so that, for each coefficient i , the expectation of w_i is zero (even conditional on any fixed value of m_i : $m_i \in \{-1, 0, 1\}$). With the modified w , the signer computes $s = f * w \pmod{q}$, and outputs the pair (m, s) as the signature of m .

Verification: To verify (m, s) , the verifier computes $t = s * h \pmod{q}$ and the deviations $Dev(s, f_0 * m)$ and $Dev(t, g_0 * m)$, as defined above in section 3, and checks that the number of deviations in each is in the range $[D_{\min}, D_{\max}]$, which is typically $[55, 87]$.

A.2 Highlighted Differences

Here is a list and brief interpretation of the changes made in R-NSS

1. For key generation, the new secret key component u replaces the public f_0 and g_0 .
2. During signing, w is calculated differently in R-NSS, using the inverse of u modulo p .
3. The stronger verification conditions in R-NSS, such as $|p^{-1}(s - m)| < B$, $|p^{-1}(t - m)| < B$, typically imply the simple deviation bounds of NSS.
4. The new distribution checks on s and t demand that they are almost normally distributed.

A.3 Suggested Parameters and Verification Steps for R-NSS

Let (N, p, q) be $(251, 3, 128)$. We take polynomials $f_1 \in \mathcal{L}(52, 52)$, $g_1 \in \mathcal{L}(36, 36)$, $u \in \mathcal{L}(88, 87)$, $m \in \mathcal{L}(80, 80)$, and $z \in \mathcal{L}(58, 58)$. However, an “efficient version” of R-NSS may set u to be a product of an element in $\mathcal{L}(7, 6)$, with one in $\mathcal{L}(7, 8)$, both f and g to be a product of two elements in $\mathcal{L}(6, 6)$, and m to be a product of two elements in $\mathcal{L}(5, 5)$, with one in $\mathcal{L}(4, 4)$. The verification procedure defines four quartiles and counts mod 3 deviations:

Quartile ranges: $I_1 = (-64, -32]$, $I_2 = (-32, 0]$, $I_3 = (0, 32]$, $I_4 = (32, 64]$.

$Dev_{1,s} = \#\{j : s_j \in I_4 | s_j - m_j = 1 \pmod{p}\} + \#\{j : s_j \in I_1 | s_j - m_j = -1 \pmod{p}\}$.

$Dev_{2,s} = \#\{j : s_j \in I_3 | s_j - m_j = 1 \pmod{p}\} + \#\{j : s_j \in I_2 | s_j - m_j = -1 \pmod{p}\}$.

Similarly for t , and set $Devsum_1 = Dev_{1,s} + Dev_{1,t}$, and $Devsum_2 = Dev_{2,s} + Dev_{2,t}$. With this notation we state the following:

Condensed Verification Criteria

1. Require $Devsum_1 \leq 10$.
2. Require $Devsum_2 \leq 18$.
3. Require $|(s', t')| < 485$.
4. Require $|s'| < 360$.
5. Require $|t'| < 360$.
6. Require between 95 and 153 coefficients of s and t to be in quartile 1.
7. Require between 50 and 100 coefficients of s and t to be in quartile 2.
8. Require between 7 and 42 coefficients of s and t to be in quartile 3.
9. Require between 0 and 14 coefficients of s and t to be in quartile 4.

B Statistical and Forgery Attacks

Here, we briefly summarize some attacks on the original version of NSS, and highlight how improvements in R-NSS defeat them. Admittedly, these attacks, by Gentry, Jonsson, Stern, and Szydlo [9], exploited specific weak features of NSS rather than its core, allowing NSS to be, in effect, “patched.” Refer to Appendix A for details on the previous version of NSS, and its differences from R-NSS.

The **Forgery Attack** provided a direct method of finding false signatures that satisfied the chief verification criterion, namely: $Dev(s, f_0 * m) < 87$, and $Dev(t, g_0 * m) < 87$. The task was to find a pair of related polynomials (s, t) that simultaneously satisfy the deviation requirements, as well as the congruence $t = s * h \pmod{q}$. Since s and t have $2N$ coefficients altogether, and the equation $t = s * h \pmod{q}$ imposes N linear constraints, there remain N degrees of freedom with which to choose the coefficients of s and t . So setting $s_i \equiv (f_0 * m)_i \pmod{p}$ and $t_j \equiv (g_0 * m)_j \pmod{p}$ for any $\lfloor N/2 \rfloor$ coefficients of s and $\lceil N/2 \rceil$ coefficients of t (i.e., about half of each), the remaining halves of s and t are determined by $g = h * f$, essentially left to chance. But the chosen half of s (resp. t) has no deviations, and the remaining half will probabilistically deviate in about $\frac{2}{3}$ of the positions, overall about $\frac{1}{3}$ of the coefficients of s (resp. t) will deviate. Since $\frac{1}{3}N \approx 84 \leq D_{\max} = 87$, this technique will generate a valid forgery after only a few iterations. See [9] for some improvements using lattice reduction to improve the chances, and “quality” of the signature.

The **Transcript Attack** recovered the private key from a long transcript (m_i, s_i) of signatures using correlations between these signatures and the private key. For a very early version of NSS, the masking polynomials w_1 , and w_2 , were chosen *at random*. Taking some m_i 's so that the average m_i was about 1, then the average s_i would be about f , the private key. About a thousand signatures were needed. Now, conceptually, we can think of that attack as a blatant example of how the coefficients f_k were jointly statistically dependent on the s_i , and m_j . Choosing w_2 so that the average above was always 0 fixed this particular flaw, but the coefficients of s still basically depended on the private key f , the message m and the polynomials w_1 and w_2 (which are not completely independent of m).

So, the same distribution concept holds: the coefficients of f depend on s in a way which is not independent of m . We can see this directly by unraveling the convolution arithmetic.

$$s_{i_0} = \sum_{j+k=i_0} f_k(m_j + v_{1,j} + pv_{2,j}).$$

Now if we only take messages, say, $m_{j_0}=1$, then we can obtain the coefficient f_k , by looking at the s_i for $i_0 = j_0 + k \bmod N$. It is not hard to see that s_i and f_k are positively correlated when m_{j_0} is always 1! Rather than attempting averaging, it is more efficient to look at the whole distribution of s_i on this transcript (which depends on f_k). With enough signatures one can tell which distribution the samples are coming from: There are only 3 shapes: the ones for $f_k = -3, 0$, or 3 . After between 10,000 and 30,000 signatures one can usually guess correctly. In this way, a transcript reveals the private key. See [9] for further details and statistics on this approach.

The scheme was revised, giving us R-NSS. One might loosely say that the attacks addressed the lack of soundness and zero knowledge inherent in NSS. Concretely, the revised version would seek verification criteria which would better bind the private key to the message, while leaking much less information. Moment balancing techniques were suggested in some of the technical notes, e.g:[14], but the main clever idea that revived NSS was the introduction of the new private key u component of f , as we introduced above when defining R-NSS.

With u , the modified scheme creates w such that the adversary knows a-priori nothing about $w \bmod p$, where before it was nearly equal to m . This effectively removed the dangerous correlation in the transcript attack[18]. Additionally, by using $u^{-1} \pmod{p}$ to construct w , with the new signature scheme, knowledge of f could be used to make $(s - m)/p \bmod q$ very small. These types of norm conditions are much stronger than the simple deviation criteria, are at least related to hard lattice problems [18], but most important at the time, obviated both of the attacks that appear in Asiacrypt '01.

Perhaps modified forgery attacks or more subtle transcript attacks would work, but this paper addresses the more fundamental issues of R-NSS, namely the mod- q reduction, and symmetry of the associated lattices.

C Cyclotomic Integers: Polynomials, Ideals, and Lattices

In this section we present some more of the background relevant to the fundamental algebraic structure behind NTRU and NSS: The ring $R = \mathbb{Z}[X]/(X^N - 1)$. We refer the reader to any standard number theory text for further background. First, we notice that the polynomial $(X^N - 1)$ factors into $X - 1$ and $X^{N-1} + X^{N-2} + \dots + X + 1$, the latter of which is irreducible[25] when N is prime and called the *Nth cyclotomic polynomial*. This induces the ring decomposition

$$R \mapsto \mathbb{Z} \times \mathbb{Z}[X]/(X^{N-1} + \dots + 1). \quad (32)$$

Projection onto the first factor is evaluation of the polynomial at 1, and for NSS, plays a smaller role in the security. The second factor is also written $\mathbb{Z}(\zeta_N)$, where (ζ_N) is a primitive N th root of unity – i.e, a solution to the cyclotomic polynomial. This very well known ring is important, because it the ring of integers in the *cyclotomic field* $\mathbb{Q}(\zeta_N)$, which is a field extension of degree $N - 1$. Both the structure of $\mathbb{Z}(\zeta_N)$ and $\mathbb{Q}(\zeta_N)$ are important for NTRU and NSS.

Lattices and Ideals: First, we recall the close relationship of polynomial rings with lattices: If A is any ring of the form $\mathbb{Z}[X]/(G(X))$ for some monic polynomial G of degree n , then the elements of A may be represented as n -dimensional vectors. The elements of A form an n -dimensional lattice isomorphic to \mathbb{Z}^n . Any sublattice S of A corresponds to an additive subgroup of A , naturally a \mathbb{Z} -module. A sublattice I of A which also has the property that for $v \in I$ and $r \in A$, the product $v * r \in I$, corresponds to an ideal of A . For NSS, when dealing with unreduced signatures, this ring is taken to be R , or $\mathbb{Z}(\zeta_N)$, and most of the lattices are principal ideals – that is, they correspond to ideals (f) and lattices generated by M_f .

The ring of integers $\mathbb{Z}(\zeta_N)$ form a Dedekind Domain, and have some properties in common with the usual integers \mathbb{Z} . The concepts of factorization and greatest common divisors must be expressed in terms of ideals in $\mathbb{Z}(\zeta_N)$. There is a unique decomposition of *ideals* into products of powers of *prime* ideals in this ring, although there are many units in the ring. While ± 1 are the only units in \mathbb{Z} , $\mathbb{Z}(\zeta_N)$ also contains the N th roots of unity and infinitely many real units. These real units may be mapped to a lattice in a $\frac{N-3}{2}$ dimensional vector space, generated by *fundamental units*.

Factorization: Not all ideals in the ring are principal ideals like (f) . The *Divisor class group* measures the fraction of ideals which are principal. For $N = 251$ this divisor class group is very large [25] and for the rational integers \mathbb{Z} it is trivial (thus having prime factorization into *elements*; one need not consider ideals). The related problem of finding the units and class number is a computationally difficult lattice problem. The factorization problem of elements, $ab \in \mathbb{Z}(\zeta_N)$ can be accomplished with prime factorization algorithms[2], but the best known algorithms use a lattice with which has as many rows as the size of the class group. For this reason, the large class group of $\mathbb{Z}(\zeta_N)$ is cited as an obstruction to factorization, and to some attacks on NSS and NTRU.

The unit group can be thought of as measuring the extent to which the map $f \mapsto (f)$ not unique, for (f) is equal to (uf) for any unit u . This contrasts with the fact that the maps $f \mapsto f^2$, and $f \mapsto f * \bar{f}$ only have kernel ± 1 , and the $2N$ th roots of 1, respectively. Sometimes a lattice basis reduction algorithm such as LLL may be able to find a generator of (f) , especially when f itself has smaller norm than most vectors in (f) . For NSS, we can think of knowledge of a polynomial $f * \bar{f}$ as the specification of a generator of (f) up to a $2N$ th root of unity – i.e, up to sign and rotation.

Given the prime factorization of ideals in $\mathbb{Z}(\zeta_N)$, the greatest common divisor (GCD) is also only defined up to a unit. So, given $a, b, f \in K$ such that a and b

have no ideal factors in common, we may compute

$$\text{GCD}(af, bf) = (f), \quad (33)$$

but this will not necessarily give us the element f . However, also knowing $f * \bar{f}$ allows us to specify f up to sign and rotation.

For published algorithms for ideal intersection, addition, multiplication, and factorization of ideals, the ideal GCD, and a Euclidean algorithm for relatively prime ideals, see [2] and [3].

D Galois Theory

In this section we summarize some relevant Galois theory of $\mathbb{Q}(\zeta_N)$. Recall the decomposition of R into $\mathbb{Z} \times \mathbb{Z}(\zeta_N)$. The field correspondent to the second factor, $K = \mathbb{Q}(\zeta_N)$, the N th cyclotomic field, is a Galois extension of \mathbb{Q} with Galois group isomorphic to $\mathbb{Z}/(N-1)\mathbb{Z}$. Representing elements of $\mathbb{Q}(\zeta_N)$ as polynomials, the automorphisms in the Galois group are the $N-1$ elements $\sigma(r)$ for $r \in \{1, 2, \dots, N-1\}$ defined by the mapping $x \mapsto x^r$. In other words, $(f^{\sigma(r)})(x) = f(x^r)$.

There is a subfield of $\mathbb{Q}(\zeta_N)$ corresponding to every subgroup of the Galois group, indexed by the factors of the integer $N-1$. The largest such proper subfield of $\mathbb{Q}(\zeta_N)$ is, explicitly, $L = \mathbb{Q}(\zeta + \zeta^{-1})$. Every element in L is a real number, and K is a quadratic extension of L with $\text{Gal}(K/L)$ consisting of only two elements: the identity and ‘complex conjugation’ σ .

Conjugation in K – namely, $\sigma(-1)$ – corresponds to the reversal of a polynomial in R . Since elements of L are fixed under conjugation, they correspond to ‘palindromes’ in R – i.e., the set of polynomials a such that $a = \bar{a}$. Given that for any $k \in K$, $k * \sigma(k)$ is a real number in L , $a * \bar{a}$ is a palindrome for any $a \in R$.

E Proof of the Galois congruences

We want to compute the polynomials of the form $f^p \pmod{p}$, so we use the basic equation

$$f^p(x) \equiv (\sum f_i x^i)^p \equiv \sum f_i^p x^{i*p} \equiv \sum f_i x^{i*p} \equiv f(x^p) \pmod{p}. \quad (34)$$

Where the coefficients are expanded, the p -multiple cross terms dropped, Fermat’s little theorem applied, and lastly, the coefficients collapsed, producing the general equation

$$f^p = (f)^{\sigma(p)} \pmod{p}, \quad (35)$$

true for all primes p , where the *Galois Conjugate* $(f)^{\sigma(r)}$ of f in R , is defined by

$$(f^{\sigma(r)})(x) = f(x^r). \quad (36)$$

In the special case where $p = N$, we have $x^N = 1$ so

$$f^N = f(1) \pmod{N}. \quad (37)$$

For the case $p = 1 \pmod{N}$, suppose that p does not divide the determinant: $\delta = \text{Det}(M_f)$. Then, there is a polynomial g such that $f * g = \delta$, so considering $f * g$ now in the ring R_p , we see that f is a unit in this ring, and we can finally draw the conclusion:

$$f^{p-1} = 1 \pmod{p}. \quad (38)$$

Lastly, in this section we make a remark about the decomposition of

$$R \mapsto \mathbb{Z} \times \mathbb{Z}[X]/(X^{N-1} + \dots + 1). \quad (39)$$

The congruence results above extend to the second factor, the ring of cyclotomic integers $\mathbb{Z}(\zeta)$, namely the general equation $f^{p-1} = 1 \pmod{p}$ is to be interpreted as an equation in $\mathbb{Z}(\zeta)$, and the automorphisms $\sigma(b)$ are the natural automorphisms in the Galois Group $\text{Gal}(\mathbb{Z}(\zeta))/\mathbb{Z}$. Equation 38 for $p = 1 \pmod{N}$ requires the assumption that p does not divide the discriminant, rather than the determinant, and equation 37 $f^N = f(1) \pmod{N}$ holds as is.

This observation can be useful, because $\mathbb{Z}(\zeta)$ is an integral domain, while R is not. The lattice arguments would hold if the lattice L is replaced with the lattice obtained by adjoining the vector 1^N of N 1's, and restricting to the sublattice where the last column is forced to have a zero coefficient. This effectively reduces all polynomials by the polynomial $X^{N-1} + \dots + 1$, and thereby removing the need to check polynomial evaluation at 1, and generally eliminating the zero divisor complications of our ring R .

F Orthogonal Congruence Attack is Polynomial-Time

In this appendix, we give a more algorithmic treatment of our orthogonal congruence attack, listing these algorithms in pseudocode, and sketch how its time-complexity is polynomial in N and the bit-length of f . As a preliminary matter, we assume that the bit-length of our input basis B_f for (f) is polynomial in N and the bit-length of f . This is a reasonable assumption for R-NSS, because when we use LLL to get B_f from $\{f * w_1, \dots, f * w_k\}$, the bit-length of B_f will be bounded polynomially by N and the bit-lengths of the $f * w_i$'s, where k will usually be a very small number (e.g., 2).

We describe the first algorithm of the attack below. Note that this algorithm will compute $f^{2^r} * \overline{v_r} \pmod{p}$ instead of $f^{p-1} * \overline{v_r} \pmod{p}$ for $r = \lfloor \log_2 p \rfloor$, with the understanding that computing the latter requires only minor adjustments based on the binary representation of $p - 1$. Let $v_0 = f$.

Algorithm 1

Input: $v_0 * \overline{v_0}$; basis B_0 for (v_0) ; prime $p \equiv 1 \pmod{N}$ with $r = \lfloor \log_2 p \rfloor = (N + 1)/2$.

Output: $\{v_0^2 * \overline{v_1}, \dots, v_{r-1}^2 * \overline{v_r}\}$ and $\{v_0 * \overline{v_0}, \dots, v_{r-1} * \overline{v_{r-1}}\}$ with $\|v_i\| < 2^{(N-1)/2}$; $v_0^{2^r} * \overline{v_r} \pmod{p}$.

1. Set $r' := 0$.
2. While $r' < r$ do
 - (a) Use $B_{r'}$ to construct a set G of vectors that generates $(v_{r'}^2)$. Let $G = H \cdot M_{v_{r'}^2}$.
 - (b) Compute the Gram matrix $H \cdot H^T = G \cdot M_{v_{r'}^2}^{-2} \cdot G^T$.
 - (c) Reduce the Gram matrix $H \cdot H^T$ to get $(A \cdot H) \cdot (A \cdot H)^T$ for known A .
 - (d) Compute the basis $(A \cdot H) \cdot M_{v_{r'}^2} = A \cdot G$ of $(v_{r'}^2)$. Note: the rows of $A \cdot H$ are short.
 - (e) Pick, say, the j th row of $A \cdot G$, call it $v_{r'}^2 * \overline{v_{r'+1}}$. Note: $\overline{v_{r'+1}}$ is the j th row of $A \cdot H$.
 - (f) Output $v_{r'}^2 * \overline{v_{r'+1}}$.
 - (g) If $r' + 1 \neq r$
 - Compute $v_{r'+1} * \overline{v_{r'+1}} = (v_{r'}^2 * \overline{v_{r'+1}}) * \overline{(v_{r'}^2 * \overline{v_{r'+1}})} * (v_{r'} * \overline{v_{r'}})^{-2}$.
 - Output $v_{r'+1} * \overline{v_{r'+1}}$.
 - Compute $(A \cdot H) \cdot M_{v_{r'+1}} = A \cdot G \cdot M_{v_{r'}^2}^{-2} \cdot M_{\overline{v_{r'+1}} * v_{r'+1}}$. This is our basis $B_{r'+1}$ of $(v_{r'+1})$.
 - (h) Increment r' .
3. Set $r' := 1$; Set $y := v_0^2 * \overline{v_1} \pmod{p}$.
4. While $r' < r$ do
 - (a) Compute $y := y^2 \pmod{p}$.
 - (b) Compute $y := y * (v_{r'} * \overline{v_{r'}})^{-2} * (v_{r'}^2 * \overline{v_{r'+1}}) \pmod{p}$. At this point, $y = v_0^{2^{r'+1}} * \overline{v_{r'+1}} \pmod{p}$.
 - (c) Increment r' .
5. Output y .

Remarks: If we had computed $v_0^{p-1} * \overline{v_r} \pmod{p}$ instead of $v_0^{2^r} * \overline{v_r} \pmod{p}$, we could have used the fact that $v_0^{p-1} \equiv 1 \pmod{p}$ to recover $\overline{v_r} \pmod{p}$, which would give us $\overline{v_r}$ exactly since $\|v_i\| < p/2$. We will have problems in step 4b if any of the v_i are zero divisors in R_p , but we can avoid this problem in step 2e by choosing a row whose corresponding lattice ideal has a determinant that is not divisible by p .¹⁰ (We know that this is the case for at least one row, since $A \cdot H$ defines the lattice \mathbb{Z}^N .)

It should be clear that Algorithm 1 is polynomial-time with the possible exception of steps 2a and 2c. Regarding step 2a, the failsafe way to get a generating set is to compute the ideals as described in section 7.4, in which the set will consist of the N^2 products $b_i * b_j$ for $b_i, b_j \in B_{r'}$. Computing these products is obviously polynomial-time, but the fairly large size of this set (N^2) will slow down step 2c. However, step 2c will still be polynomial-time, since both H 's dimension and the bit-lengths of H 's entries will be bounded polynomially in N .

¹⁰ See Appendix E for more on how zero divisorship relates to determinants of lattice ideals.

Although this completes the proof that Algorithm 1 is polynomial-time, there are also some practical considerations to address.

In practice, we have found that small subsets (on the order of N) of the N^2 products above are often generating sets. For example, if $b_i = a_i * v_{r'}$ and $b_j = a_j * v_{r'}$ where a_i and a_j are relatively prime, which will often be the case, the a_i^2 and a_j^2 will also be relatively prime, and we can generate $(v_{r'}^2)$ from the $2N$ row vectors of b_i^2 and b_j^2 .

Another practical issue is that LLL usually finds vectors much shorter than its worst-case guarantee – meaning that the v_i 's in Algorithm 1 will usually have norms much less than $2^{(N-1)/2}$. This helps us in two ways. First, each lattice reduction step will be faster, because the input bases will already be substantially reduced. Second, fewer lattice reductions will be necessary, because we can choose p smaller. In the future, we anticipate performing computer simulations to determine the actual time it takes to recover an R-NSS private key. Below, we give two more algorithms, for which we omit specific proofs of time-complexity.

Algorithm 2

Input: $\{v_0^2 * \overline{v_1}, \dots, v_{r-1}^2 * \overline{v_r}\}$ and $\{v_0 * \overline{v_0}, \dots, v_{r-1} * \overline{v_{r-1}}\}$ from above; $\overline{v_r}$ exactly; p above also satisfies $\gcd(p-1, 2N-1) = 1$.

Output: v_0^{2N} .

1. Let $P_1 > 2\|v_0^{2N}\|$ be a prime number satisfying $P_1 \equiv -(2N-1) \pmod{p-1}$ – i.e., $P_1 = (p-1)c_1 - 2N + 1$.
2. Compute $v_0^{p-1} * \overline{v_r} \pmod{P_1}$ from the sequences $\{v_0^2 * \overline{v_1}, \dots, v_{r-1}^2 * \overline{v_r}\}$ and $\{v_0 * \overline{v_0}, \dots, v_{r-1} * \overline{v_{r-1}}\}$ generated in Algorithm 1.
3. From $v_0^{p-1} * v_r \pmod{P_1}$ and $\overline{v_r}$, compute $v_0^{p-1} \pmod{P_1}$.
4. Compute $(v_0^{p-1})^{c_1} \equiv v_0^{P_1 + (2N-1)} \equiv v_0^{2N} \pmod{P_1}$. Since $P_1 > 2\|v_0^{2N}\|$, $v_0^{2N} \pmod{P_1}$ gives v_0^{2N} exactly.

Algorithm 3

Input: f^{2N} .

Output: f , up to rotation and sign.

1. Compute the smallest positive odd integer b such that $N-b$ is a primitive root modulo N . (For example, $b = 3$ for $N = 251$ since 248 is a primitive root modulo 251.)
2. Let $P_2 > 2\|f^b(X) * f(X^{N-b})\|$ be a prime number satisfying $P_2 \equiv -b \pmod{N}$ – i.e., $P_2 = 2c_2N - b$.
3. Compute $(f^{2N})^{c_2} \equiv f^b * f^{P_2} \equiv f^b(X) * f(X^{N-b}) \pmod{P_2}$. Since $P_2 > 2\|f^b(X) * f(X^{N-b})\|$, we get $f^b(X) * f(X^{N-b})$ exactly.
4. Let α be a complex N th root of unity. Compute $f^{2N}(\alpha)$.
5. Pick a $2N$ th root of $f^{2N}(\alpha)$ and call it $f(\alpha)$.
6. Iteratively compute $f(\alpha^{(N-b)^{r+1} \pmod{N}})$ from $f(\alpha^{(N-b)^r \pmod{N}})$ using the polynomial $f^b(X) * f(X^{N-b})$.
7. Use the $N-1$ values of $f(\alpha^{(N-b)^r \pmod{N}})$, as well as $f(1)$, to compute f .