

Characterizing overstretched NTRU attacks

Gabrielle De Micheli, Nadia Heninger and Barak Shani
University of Pennsylvania

Abstract

Overstretched NTRU, an NTRU variant with a large modulus, has been used as a building block for several cryptographic schemes in recent years. Recently, two lattice *subfield attacks* and a *subring attack* were proposed that broke some suggested parameters for overstretched NTRU. These attacks work by decreasing the dimension of the lattice to be reduced, which improves the performance of the lattice basis reduction algorithm. However, there are a number of conflicting claims in the literature over which of these attacks has the best performance. These claims are typically based on experiments more than analysis. Furthermore, the metric for comparison has been unclear in some prior work. In this paper, we argue that the correct metric should be the lattice dimension. We show both analytically and experimentally that the subring attack succeeds on a smaller dimension lattice than the subfield attack for the same problem parameters, and also succeeds with a smaller modulus when the lattice dimension is fixed.

Keywords: overstretched NTRU, subfield attack, subring attack

1 Introduction

NTRU is a public key cryptosystem introduced by Hoffstein, Pipher and Silverman [HPS96, HPS98]. NTRU serves as a basis for many cryptographic protocols (e.g. [HHGP⁺03, LATV12, GGH13]) and is believed to remain secure in the presence of quantum computers. See the survey [Ste14] for a complete description of the NTRU cryptosystem and its applications.

In this paper, we consider NTRU in the cyclotomic field $K = \mathbb{Q}[x]/(x^n + 1)$, with n a power of 2. Let $R = \mathbb{Z}[x]/(x^n + 1)$ be the ring of integers in K , and let $R_q = \mathbb{Z}_q[x]/(x^n + 1)$ for some integer q . The private key consists of two polynomials $\mathbf{f}, \mathbf{g} \in R_q$, with \mathbf{f} invertible, and the public key \mathbf{h} is defined by $\mathbf{h} := \mathbf{g}\mathbf{f}^{-1}$. The coefficients of \mathbf{f} and \mathbf{g} are chosen to be small, and follow a given distribution. The uniform distribution over $\{-1, 0, 1\}$ is most common in practice. When the coefficients of \mathbf{f} and \mathbf{g} are drawn from a Gaussian distribution with standard deviation $\sigma \approx \sqrt{q}$, the distribution of the public key is statistically indistinguishable from a uniformly sampled element [SS11]. The NTRU key recovery problem is to recover the private key (\mathbf{f}, \mathbf{g}) given $\mathbf{h} = \mathbf{g}/\mathbf{f} \in R_q$.

We focus on a variant of NTRU, called *overstretched NTRU*, which sets the size of the modulus q to be super-polynomial in n . The main motivation for choosing such a large modulus q is that it allows several arithmetic operations on small elements before the result “wraps around” mod q , thus resulting in a somewhat homomorphic encryption system [LATV12, BLLN13]. In overstretched NTRU, the scheme is insecure if an attacker can find any pair of polynomials $\mathbf{a}, \mathbf{b} \in R_q$ with small

coefficients such that $\mathbf{h} = \mathbf{b}\mathbf{a}^{-1}$, and so we relax key recovery to the problem of finding any such polynomials $\mathbf{a}, \mathbf{b} \in R_q$.¹

Cryptanalyses of NTRU A lattice attack on NTRU was first given by Coppersmith and Shamir [CS97]. This attack takes place over the full field K , and the resulting lattice is of dimension $2n$. The lattice dimension affects the performance of lattice basis reduction algorithms in two ways. First, the approximation factor achieved by short vectors grows with the dimension, so we would expect a higher-dimension lattice to reduce to a “worse” basis than a lower-dimension lattice. Second, the running time of basis reduction algorithms depends on the dimension. It is desirable for both of these reasons to develop attacks that use lattices of smaller dimension. May [May99] slightly modified the Coppersmith-Shamir lattice attack using *projection*; however the resulting dimension is not significantly lower than the original dimension.

Overstretched NTRU gives attackers greater flexibility in the lattice construction, as was first pointed out by Gentry and Szydlo [GS02], and rigorously developed in the following works. The first cryptanalysis of overstretched NTRU was given independently by Albrecht et al. [ABD16] and Cheon et al. [CJL16]. These works present *subfield attacks* that exploit the presence of subfields in K . The main result presented in these works is the fact that if q is chosen to be exponentially large in the security parameter, while n is polynomial in it, then the attack runs in polynomial time. This is due to the fact the resulting lattice dimension is significantly smaller than n .

Kirchner and Fouque [KF17] developed a new attack against overstretched NTRU called a *subring attack*. The subring attack allows more flexibility in choosing different (larger) dimensions that cannot be achieved in the subfield attacks. In addition, using the Pataki-Tural lemma [PT08], they proved a very strong result showing that, despite the larger dimension, the “full field” attack does not perform worse than the subfield and subring attacks. Roughly speaking, the result shows that lattice reduction algorithms in the full field attack already exploit the existing subfields. From the complexity point of view, this means that asymptotically, both attacks run in polynomial time for the same set of parameters. We remark that this result relies on some technical conditions that are not known to hold in general.

From the asymptotic point of view, if one accepts these technical conditions, there is no benefit in using the new attacks. On the other hand, since the full field attack is currently impractical for modest parameter sizes (e.g. for $n > 2^{12}$), it remains of interest to study dimension reduction attacks. Our experiments, as well as those of previous works, show that the running time of the attack in subfields is significantly shorter than in the full field. In [KF17], the authors experimentally compare the subring attack to similar experiments for the subfield attacks from [ABD16], and conclude from experiments that the subring attack “performs better” than the subfield attack if one wishes to minimize the ratio between n and q . In a subsequent work, Duong et al. [DYT17] experimentally compare different subfields and observe that the subfield attack of [ABD16] sometimes “performs better” than, and sometimes is equivalent to, the subring attack. These comparisons are mainly experimental and not analytical.

A deeper look into these comparisons shows that they happen between lattices of different

¹This is also true for NTRU. However, the small modulus restricts the number of such polynomials.

dimensions, where attacks on larger-dimension lattices seem to produce better results, that is, that the attack succeeds on a smaller modulus q for a given degree n . We give concrete examples in the end of Section 3.2. The benefit of this comparison is questionable, as it follows from [KF17] that among these lattice attacks, the full field, which has the largest dimension, is expected to achieve the lowest ratio between n and q . Furthermore, we remark that the point of the subfield and subring attacks is to decrease the dimension, so increasing the dimension is in opposition to the goal of the attack construction. It should be noted that in general the lattice dimension is not the only parameter that affects the running time or approximation factor of lattice basis reduction algorithms. However, in all of our experiments, as in the reported experiments of previous works, decreasing the lattice dimension reduces the running time (see Table 3).

Our results The main goal of our work is to analyze the relative performance of the different lattice attacks and thus resolve the conflicting claims in previous work. In contrast to prior work, we focus on analysis and use experiments to validate our analysis. Our analysis focuses on the lattice dimension, following our claim that this is the correct metric for comparison, as we explain above.

Our contributions include:

1. We give formal justification for the projection technique of May [May99] and May and Silverman [MS01], which is key to the subring attack. We formalize the condition under which this technique works, and explain its relation with some standard assumptions on NTRU. This analysis is missing in previous work, and thus we fill a theoretical gap in the subring attack.
2. We show that the subring attack is expected to perform better than the subfield attack. This result resolves the incompatible claims in previous works. In short, for fixed parameters n and q , it is possible to obtain a lower dimension lattice in the subring attack by discarding more equations from the lattice using the projection technique. As a consequence, for a given degree n and fixed dimension lattice, the subring attack is expected to succeed for a smaller modulus q .

We summarize our results in the following informal theorem statements.

Main Result (Informal). *Consider the NTRU problem with polynomial degree n and modulus q . Let L be a subfield of K such that $[K : L] = r$, and let $n' = n/r$. Then, for sufficiently large n , we have the following.*

1. *Consider the subfield and subring lattice constructions in dimension $2n'$. The subring lattice is expected to contain shorter vectors than the subfield lattice, and thus can also solve the NTRU problem with smaller modulus. The ratio between the Euclidean norm of the desired lattice vectors in the subfield and subring lattices is approximately $\sqrt{2r/(r+1)}$. If these are the shortest vectors, then the ratio between the feasible moduli for which each attack works approaches 2 as r increases.*
2. *Furthermore, if we use the projection technique to decrease the lattice dimension below $2n'$, the subring attack is expected to solve the NTRU problem on a smaller dimension lattice than the*

subfield attack, and using a smaller block size for the BKZ algorithm, by finding a non-zero integral multiple of the desired vectors.

Our analysis does not show that these desired vectors are the shortest vectors in the corresponding lattices. Thus, our bounds may not be tight. We present experimental evidence suggesting that these bounds are conservative.

Our result focuses on the structure of the lattices more than actual implementations of the attacks. In particular, we fix the subfield index and analyze the asymptotics of the length of short vectors in the lattices. Implementations of the attacks would try to optimize the choice of the subfield with respect to the degree of the field. An analysis of such an optimization seems challenging.

2 Preliminaries

We introduce definitions and results from algebraic number theory (see [Sam70] for more details) and background on lattices. Throughout the paper, we distinguish between representations of elements and denote the ring representation by $\mathbf{a} \in R$, and the vector representation by $a \in \mathbb{Z}_q^n$. The Euclidean norm $\|\mathbf{f}\|$ is taken to be the norm of the corresponding vector consisting of the polynomial coefficients, i.e. $\|f\|$. Also, we use the notation $[x]_q$ to indicate that x is taken mod q .

2.1 Background in number theory and lattices

Number fields A *number field* is a finite field extension of \mathbb{Q} . Its degree is $[K : \mathbb{Q}]$. The ring of integers \mathcal{O} of a number field K is the set of algebraic integers contained in K . For any field K and subfield L of K , we define $r = [K : L]$ to be the index of the subfield we consider. If $n' = [L : \mathbb{Q}]$ and $n = [K : \mathbb{Q}]$, then $r = n/n'$.

Cyclotomic fields Let ζ_m be a primitive m^{th} root of unity. We define the m^{th} *cyclotomic field* to be $K = \mathbb{Q}(\zeta_m)$. In this paper, as in most applications, we are interested in the case where m is a power of two. Then, the m^{th} cyclotomic polynomial is $x^n + 1$, where $n = \phi(m)$ and ϕ is Euler's phi function, and it holds that $K \cong \mathbb{Q}[x]/(x^n + 1)$.

Relative norm and trace Let K be a number field and L a subfield of K . For any element $a \in K$, we consider the map $m_a : x \mapsto ax$, for $x \in L$. The trace of $a \in K$, denoted $Tr_{K/L}(a)$, is the trace of m_a , and the relative norm of $a \in K$, denoted $N_{K/L}(a)$, is the determinant of m_a . The trace is additive, $Tr(x + x') = Tr(x) + Tr(x')$, and the norm is multiplicative, $N(xx') = N(x)N(x')$ for any $x, x' \in K$. More specifically, if K is a Galois extension of \mathbb{Q} and we define $G = Gal(K/L)$, we have

$$Tr_{K/L}(a) = \sum_{\sigma \in G} \sigma(a) \quad \text{and} \quad N_{K/L}(a) = \prod_{\sigma \in G} \sigma(a).$$

The embeddings $\sigma \in G$ permute or conjugate the coefficients of $x \in K$ in the canonical embedding representation. Hence, we have that $\forall \sigma \in G, \|\sigma(x)\| = \|x\|$.

In the paper, we will consider a number field K and a subfield L of index r . We generally write $N_{K/L} = \prod_{\sigma \in G} \sigma$ for the relative norm, as above. However, we sometimes enumerate the embeddings. Specifically, as one of the embeddings is the identity function, we set $\sigma_1 = \text{Id}$. To prevent confusion, we clarify that while the canonical embeddings are used, the norms are taken with respect to the coefficients.

Lattices A *lattice* is a discrete additive subgroup of \mathbb{R}^n . For any full-rank lattice L of dimension n there exists a basis $B = \{b_1, b_2, \dots, b_n\}$, consisting of linearly independent vectors b_i . When $n \geq 2$, there are infinitely many such bases. We will represent an n -dimensional lattice as an $n \times n$ matrix where the rows are given by the basis vectors b_i , and write $\mathcal{L}(B)$ for the lattice generated by basis matrix B . We denote by λ_1 the shortest non-zero vector of the lattice, i.e., $\lambda_1(L) = \min_{v \in L \setminus \{0\}} \|v\|$. More generally, we write λ_i to denote the i^{th} successive minimum of the lattice.

3 Characterization of NTRU attacks

In this section we present the lattice attacks on NTRU and overstretched NTRU. We begin with an attack on NTRU in the full field, and present ideas to reduce the dimension of the resulting lattice. We find it helpful to view the new attacks on overstretched NTRU as a modification of this lattice in light of these ideas.

3.1 The full field attack

We first present the NTRU lattice in the “full field”, which is the attack of Coppersmith and Shamir [CS97]. Consider the equation

$$\mathbf{f} \cdot \mathbf{h} = \mathbf{g}, \tag{1}$$

which relates the private and public keys in NTRU. It is well known (see for example [LN97, Section 2.3]) that each coefficient g_j of the polynomial \mathbf{g} is an n -dimensional linear function in the coefficients f_i of the polynomial \mathbf{f} . One can therefore represent Equation (1) as a system of n linear equations in the unknowns f_i, g_i . This linearity, along with the fact that the coefficients of \mathbf{f} and \mathbf{g} are relatively small, suggests using the structure of a lattice when trying to solve the NTRU problem. This was noted by Coppersmith and Shamir [CS97], who considered the lattice generated by the rows of the following $2n \times 2n$ matrix

$$A_{full} = \begin{pmatrix} \mathbf{I}_n & \mathcal{M}_{\mathbf{h}} \\ 0 & q\mathbf{I}_n \end{pmatrix},$$

where \mathbf{I}_n is the n -dimensional identity matrix and $\mathcal{M}_{\mathbf{h}}$ represents multiplication in R by the public key \mathbf{h} . The lattice generated by A_{full} , which we call $\mathcal{L}(A_{full})$, contains the vector $(f, g) = (f_0, \dots, f_{n-1}, g_0, \dots, g_{n-1})$, because $(f_0, \dots, f_{n-1})\mathcal{M}_{\mathbf{h}} \pmod{q} \equiv (g_0, \dots, g_{n-1})$. In the following descriptions, for any lattice vector, we use the term “ f part” for the part of the vector corresponding to the identity matrix, and “ g part” for the part corresponding to multiplication by \mathbf{h} (or by $N_{K/L}(\mathbf{h})$ in the subfield lattice).

The vector (f, g) is a *short vector* (i.e. has relatively small Euclidean norm), and is most likely a vector of the smallest non-zero length in the lattice $\mathcal{L}(A_{full})$. In addition, Coppersmith and

Shamir [CS97] note that one can derive useful information to recover the secret key even when a multiple of (\mathbf{f}, \mathbf{g}) is found, for multiples of relatively small norm (yet larger than (\mathbf{f}, \mathbf{g})). See [CS97] for more details.

Dimension reduction methods It follows that the NTRU problem can be reduced to computing short vectors in the lattice $\mathcal{L}(A_{full})$. The main obstacle to this approach is the lattice dimension, which is too large to run a practical attack for realistic parameters. May, in an unpublished work [May99], and later with Silverman [MS01] describe some methods to reduce the dimension of a certain class of lattices. Their work can be applied to the NTRU lattice. We give a short description of these methods.

First, we observe that Equation (1) holds when the polynomials \mathbf{f} and \mathbf{g} are replaced by $\mathbf{f}x^i, \mathbf{g}x^i$, as $\mathbf{f}x^i \cdot \mathbf{h} = \mathbf{g}x^i$. In NTRU, multiplying polynomials by x^i rotates the polynomial’s coefficients in the originally suggested ring $\mathbb{Z}[x]/(x^n - 1)$. The vector $\mathbf{f}x^i$ is called a rotation, or *shift*, of \mathbf{f} . In the ring $\mathbb{Z}[x]/(x^n + 1)$ each shift of \mathbf{f} also changes the sign of some of the coefficients. However, the Euclidean norms of \mathbf{f} and any of its shifts are equal. Hence, the NTRU lattice does not contain a unique shortest vector in general.

Zero forcing: In order to increase the ratio $\lambda_2(\mathcal{L}(A_{full}))/\lambda_1(\mathcal{L}(A_{full}))$, note that if one knows some zero coefficients in \mathbf{f} , then multiplying the corresponding rows in the matrix block \mathbf{I}_n by a sufficiently large scaling factor θ (that is, if $f_j = 0$, we multiply the j -th row of \mathbf{I}_n by θ) will “eliminate” the shifts whose j -th coordinate is non-zero from the resulting lattice. In [MS01, Section 7], the authors claim that this technique, called *zero forcing*, helps reduce the “effective dimension” of the lattice $\mathcal{L}(A_{full})$ by a factor proportional to the number of known zeros. To increase the ratio $\lambda_2(\mathcal{L}(A_{full}))/\lambda_1(\mathcal{L}(A_{full}))$, one has to eliminate all the shifts from $\mathcal{L}(A_{full})$. Note that instead of multiplying the rows by θ , one can remove these rows from $\mathcal{L}(A_{full})$ to achieve an actual dimension reduction. However, this does not help in eliminating all shifts.

Projection: When the polynomials \mathbf{f} and \mathbf{g} have coefficients chosen from $\{-1, 0, 1\}$, or other sets with small values compared to q , the system of n equations coming from (1) is likely to be overdetermined. Thus, even if some equations are discarded, it is still likely that the only solution in $\{-1, 0, 1\}$ is the original solution (and its shifts). This observation suggests that one can “throw away” some of the columns in \mathcal{M}_h (and the corresponding ones in $q\mathbf{I}_n$) to reduce the dimension of $\mathcal{L}(A_{full})$. We present a concrete analysis of this method, called *projection*, in Section 4.1.

3.2 Variants of the Coppersmith–Shamir attack

We now describe the subfield attacks of [ABD16, CJL16] and subring attack of [KF17] in light of the dimension reduction methods presented above. These attacks are variants of the Coppersmith–Shamir attack that exploit the fact that it suffices to find relatively small multiples of the private key. When the modulus q is large, as in overstretched NTRU, these multiples can be quite large as well. In this case, the ring structure aids in reducing the lattice dimension. We restrict the following discussion to n a power of two. We focus on finding a small multiple of \mathbf{f} , as once we find a multiple $\alpha\mathbf{f}$, it is straightforward to derive a small multiple of \mathbf{g} by computing $\alpha\mathbf{f}\mathbf{h} = \alpha\mathbf{g}$.

The three attacks apply the same method in reducing the f part of vectors in $\mathcal{L}(A_{full})$, but differ in their construction of the g part. We start by explaining the reduction of the f part, and then present for each attack its g part reduction.

The main observation in these attacks is the existence of elements with known zero polynomial coefficients. Focusing only on these elements allows one to modify the matrix A_{full} to have lower rank and thus it spans a lattice \mathcal{L} of lower dimension. The attacks exploit the following three properties:

- *The images under $N_{K/L}$ and $Tr_{K/L}$ are sparse polynomials in K .* This holds because elements of a proper subfield $L \subset K$, when represented as polynomials, are sparse (i.e. have few non-zero coefficients). The number of non-zero coefficients is bounded by the degree of L as a number field. This property follows from the fact that the norm $N_{K/L}$ and the trace $Tr_{K/L}$ are mappings into L .
- *$N_{K/L}(\mathbf{f})$ is a multiple of \mathbf{f} .* We have that $N_{K/L}(\mathbf{f}) = \prod_{\sigma \in G} \sigma(\mathbf{f})$. Since one of these embeddings is the identity map, we have $N_{K/L}(\mathbf{f}) = \mathbf{f} \cdot \prod_{\sigma \in G \setminus \{Id\}} \sigma(\mathbf{f})$.
- *$\|N_{K/L}(\mathbf{f})\|$ is relatively small if $\|\mathbf{f}\|$ is small.* Recall that $\|\sigma(\mathbf{f})\| = \|\mathbf{f}\|$. Thus, an upper bound on the product $\prod_{\sigma \in G} \sigma(\mathbf{f})$ can be given, for example, by iteratively applying the Cauchy–Schwarz inequality as done in [CJL16, Lemma 1]. For example, when $[K : L] = 2$ then $N_{K/L}(\mathbf{f}) = \mathbf{f}\sigma(\mathbf{f})$, and so $\|N_{K/L}(\mathbf{f})\| \leq \|\mathbf{f}\| \|\sigma(\mathbf{f})\| \sqrt{n} = \|\mathbf{f}\|^2 \sqrt{n}$.

These facts allow us to reduce the search space for multiples of \mathbf{f} to those that lie in some subfield of K . That is, the subfield attack restricts the space to multiples of $N_{K/L}(\mathbf{f})$, which are also multiples of \mathbf{f} . These multiples will not be as small as \mathbf{f} in general. However, working in a subfield allows one to greatly reduce the lattice dimension in the following way. Since $N_{K/L}(\mathbf{f}) \in L$ we know the positions of its zero coefficients. In this case, we can use the zero forcing technique and remove² the rows (and subsequently zero columns) that correspond to known zero coefficients of $N_{K/L}(\mathbf{f})$ from the matrix A_{full} . The result is a modified lattice whose f part only contains vectors that correspond to elements in the subfield L . We give an example:

Example 1. *Consider a subfield L such that $[K : L] = 2$. Then all coefficients of x^{2i+1} in $N_{K/L}(\mathbf{f})$ are zero. Instead of searching for a multiple of $(f, g) = (f_0, \dots, f_{n-1}, g_0, \dots, g_{n-1})$ in the original lattice \mathcal{L} , we restrict to vectors whose f parts are multiples of $N_{K/L}(\mathbf{f})$. The first n coordinates of such vectors are of the form $(u_0, 0, u_2, 0, \dots, 0, u_{n-2}, 0)$. Thus, one can remove every second row in the top half of A and the corresponding zero columns. This results in a lattice of dimension $3n/2$ instead of $2n$.*

For a subfield L of degree n' , this step gives the following $(n' + n) \times (n' + n)$ matrix

$$A_{subring} = \begin{pmatrix} \mathbf{I}_{n'} & \mathcal{M}'_{\mathbf{h}} \\ 0 & q\mathbf{I}_n \end{pmatrix},$$

²The zero forcing method described earlier suggests that if some zeros are known, one can multiply the corresponding rows by some value θ in order to eliminate the shifts. We remark that while some zeros are known in this approach, the symmetry in this sequence of zeros will not help in removing the shifts. Indeed, in the example above, even if one multiplies by θ instead of removing the rows, every second shift still appears in the lattice, which gives $n/2$ shifts. Moreover, these are exactly the shifts that appear in the modified lattice, i.e. when the rows are deleted.

that generates the (full) *subring* lattice, where $\mathcal{M}'_{\mathbf{h}}$ is the modified matrix³. This lattice contains the $(n + n')$ -dimensional vector

$$(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{f})\mathbf{h}) = \left(\prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f}) \cdot \mathbf{f}, \prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f}) \cdot \mathbf{g} \right).$$

As mentioned above all attacks reduce the f part in this way. We now explain how each attack modify the g part to achieve further dimension reduction.

The subfield attacks [ABD16, CJL16] We focus on the *norm* attack of [ABD16], because it performs slightly better when the polynomials are balanced (that is, the polynomials \mathbf{f} and \mathbf{g} have approximately the same number of non-zero coefficients), which is our case of interest. From the multiplicity of the norm function we have $N_{K/L}(\mathbf{f})N_{K/L}(\mathbf{h}) = N_{K/L}(\mathbf{f}\mathbf{h}) = N_{K/L}(\mathbf{g})$. The latter is sparse and of small size, and is suitable as a g part if one replaces the matrix block $\mathcal{M}_{\mathbf{h}}$ by $\mathcal{M}_{N_{K/L}(\mathbf{h})}$ in A_{full} .⁴ As above, the polynomial $N_{K/L}(\mathbf{h})$ is sparse, and so there are zero columns in $\mathcal{M}_{N_{K/L}(\mathbf{h})}$. These columns can be removed along with the corresponding rows in the bottom-right block. This result is a similar dimension reduction as above. The resulting lattice is generated by the rows of the $2n' \times 2n'$ matrix

$$A_{subfield} = \begin{pmatrix} \mathbf{I}_{n'} & \mathcal{M}'_{N_{K/L}(\mathbf{h})} \\ 0 & q\mathbf{I}_{n'} \end{pmatrix},$$

where n' is the degree of the subfield L and $\mathcal{M}'_{N_{K/L}(\mathbf{h})}$ is an $n' \times n'$ matrix representing multiplication by $N_{K/L}(\mathbf{h})$ in L . Note that this lattice contains the short vector $(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{g})) \in \mathbb{Z}_q^{2n'}$. As mentioned above, it is sufficient to find a short multiple of this vector. Finally, once a multiple of the vector

$$(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{g})) = \left(\prod_{\sigma \in G} \sigma(\mathbf{f}), \prod_{\sigma \in G} \sigma(\mathbf{g}) \right)$$

is found, then specifically a short multiple $\alpha N_{K/L}(\mathbf{f})$ of \mathbf{f} is found. We then “lift” $\alpha N_{K/L}(\mathbf{f})$ back to the full field using the canonical inclusion map of $L \subset K$, and as explained above, compute $\alpha \prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f})\mathbf{g} = \alpha N_{K/L}(\mathbf{f})\mathbf{h}$ in the full field. The resulting pair of elements in R_q is $(\alpha \prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f}) \cdot \mathbf{f}, \alpha \prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f}) \cdot \mathbf{g})$, which is a relatively small multiple of (\mathbf{f}, \mathbf{g}) .

The subring attack [KF17] The subring attack uses the projection technique on the g part in $A_{subring}$, that is removing columns of $\mathcal{M}'_{\mathbf{h}}$ in $A_{subring}$, to reduce the lattice dimension. For rigorous analysis, the columns should be chosen independently. This results in the following $(n' + m) \times (n' + m)$ matrix

$$A_{subring}^{proj} = \begin{pmatrix} \mathbf{I}_{n'} & \overline{\mathcal{M}'_{\mathbf{h}}} \\ 0 & q\mathbf{I}_m \end{pmatrix},$$

³The $n' \times n$ matrix $\mathcal{M}'_{\mathbf{h}}$ can be thought of as multiplication by \mathbf{h} in L .

⁴The *trace* attack [CJL16] replaces $\mathcal{M}_{\mathbf{h}}$ by $\mathcal{M}_{Tr_{K/L}(\mathbf{h})}$. Then the g part of vectors in the resulting lattice is of the form $N_{K/L}(\mathbf{f})Tr_{K/L}(\mathbf{h})$, which can be shown to be a relatively short vector. See [CJL16] for full details.

where $\overline{\mathcal{M}}'_h$ is the projected matrix of dimension $n' \times m$, that is, the number of removed columns is $n - m$. This lattice contains the vector $(N_{K/L}(\mathbf{f}), \overline{N_{K/L}(\mathbf{f})\mathbf{h}})$, where $\overline{N_{K/L}(\mathbf{f})\mathbf{h}}$ is the corresponding projected vector. Once a small vector in this lattice is found, the lifting is done as in the subfield attack.

For this vector to be the shortest vector in the lattice we need the system of equations coming from (1) to be overdetermined, such that after removing $n - m$ equations it is still determined. We give an analysis below on the number of columns that can be kept while the system remains fully determined, but in general one can experiment as there are at most n possibilities.

It is natural to set $m = n'$ in order to have equal dimension in both attacks. We show below that if the system is determined with n' equations in the subfield, then it is also determined with n' equations in the subring. In this case, the subring lattice is denser than the subfield lattice (because the matrix entries are smaller), and thus we expect it to contain shorter vectors. We will quantify this in Section 4.2. Note that the subring attack offers more flexibility than the subfield attack, as one can keep more than $2n'$ columns of $A_{subring}$. It is worth noting that one could keep fewer than $2n'$ columns, assuming the system of equations still remains determined, but this is not specific to the subring attack, and also applies in the subfield case.

Previous comparison of the attacks We now discuss the comparison between the attacks in previous works, and explain why it is problematic. The initial experimental results of [ABD16] for the subfield attack focused on finding the minimal modulus q for which the NTRU problem can be solved with the subfield attack on a fixed n using the LLL algorithm. The attacker has a choice of different subfields that could be used in this attack, which would result in different dimensions for the lattice. For the subring attack, [KF17] directly compared against these experiments and demonstrated that the subring attack, over the same subfield L , succeeded for a smaller choice of modulus q . In response, [DYT17] ran experiments over subfields of smaller index demonstrating recovery for a smaller modulus using the subfield attack. However, we observe that by [KF17, Theorem 9] we already know that, under some conditions, working over a subfield of smaller index (including the full field) will not give worse results despite the increase in dimension (however, see our experimental results for the full field compared to the other attacks in Table 2).

In a direct comparison of experiments, both [KF17] and [DYT17] use lattices of larger dimension compared to [ABD16]. Thus, it is not clear whether the experiments remain comparable. Table 1 presents one example, for $n = 2^{11}$, that demonstrates these comparisons. We observe that [KF17] applied the subring attack in dimension 638, but the subfield attack achieves a smaller q with dimension 512. It is possible that the former running time might be shorter, but this was not claimed by [KF17], and our experiments have shown a close relationship between the dimension and running time. More examples can be found in the previous work. Table 1 presents some additional comparisons we have made: for $n = 2^{11}$ and $\log(q) = 81$, the two attacks are both successful, in spite of the fact that it was claimed as an improvement in [KF17] for a fixed subfield index. On the other hand, for $n = 2^{12}$ and $\log(q) = 157$ our experiments could not succeed by taking a smaller-index subfield. We conclude that it is not clear from the experimental results in prior work in what sense one attack “performs better” than the other.

Attack	$\log n$	$\log q$	$\log r$	Dimension
Subfield [ABD16]	11	165	4	256
Subring [KF17]	11	115	4	638
Subfield [ABD16]	11	95	3	512
Subring [KF17]	11	81	3	856
Subfield (our exp.)	11	81	2	856
Subfield [DYT17]	11	72	2	1024
Subring [KF17]	12	157	4	686
Subring (our exp.) (failed)	12	157	3	686
Subfield (our exp.) (failed)	12	157	3	686

Table 1: **Experimental comparisons between the subfield and subring attacks, from prior work and our own experiments.** Some claimed improvements from prior work (smaller q) for the subring and subfield attacks have used larger dimension lattices. For $\log(n) = 11$, we show that the modulus q leading to a successful subring attack also leads to a successful subfield attack on a lattice of the same (or smaller) dimension. For $\log(n) = 12$ on the other hand, this is not always the case.

As explained above, the main feature of the subfield and subring attacks is that they allow the attacker to decrease the lattice dimension, bringing larger parameters into feasible range in practice. Since at the limit, by the aforementioned result of Kirchner and Fouque, the full field attack will work for a given n and q if the subfield attack does (using LLL or a fixed block size in BKZ), we claim that comparing results using lattices of larger dimension is not a fair method. We propose that a better metric is to first fix the lattice dimension, and then evaluate which attack succeeds with a smaller modulus q . (One should also compare the running times, but they appear to be similar as the lattices are very similar; for fixed parameters they have the same volume. See Table 3 for some timings of the lattice reduction of both attacks.) An alternative approach would be to first fix the modulus, and compare the smallest dimension, using the projection technique, for which each attack succeeds. We present our own experimental results for these two comparisons in Tables 3 and 2. In the following section we show analytically that the subring attack is expected to perform better than the subfield attacks under this metric.

4 Main results

Our main contribution is a full characterization of the various attacks on overstretched NTRU. In Section 4.1, we give a detailed analysis of the applicability of the projection technique to NTRU lattices. This technique is essential for the subring attack to hold, but the authors of [KF17] (as well as [May99]) rely on experimental evidence to justify its validity. Having fully proven the subring attack, we compare it to the subfield attack in Section 4.2. We show that the (projected) subring lattice is expected to contain shorter vectors than the subfield lattice, and derive applications of this result to the attacks.

4.1 The projection technique

The key to the projection technique is the assumption that the system of equations corresponding to Equation (1) is overdetermined, and therefore some of the equations can be discarded without introducing undesirable solutions. The aim of this section is to formalize this assumption, explain its relation to standard assumptions on NTRU and derive concrete results using it. Some of these details are missing in the results presented in [KF17, Section 3.2]. Indeed, they simply assume, backed up by experimental evidence, that a sufficiently short vector in the lattice must be a short multiple of the key.

In the following, the matrix A is any of the matrices in the attacks above and $\mathcal{L}(A)$ is the lattice generated by A . Let \mathcal{L}' be the projected matrix, as explained in Section 3.2, of dimension $n' + d$. The matrix \mathcal{M} is the corresponding matrix at the upper right quadrant of A .

We introduce the notion of *discrepancy* [KN74, DT97], which measures how close the distribution of a sequence of points is to being equidistributed. Formally, the discrepancy $\mathcal{D}(\Gamma)$ of a sequence of points $\Gamma = \{\gamma_1, \dots, \gamma_n\}$ in the interval $[0, 1]$ is defined as

$$\mathcal{D}(\Gamma) = \sup_{J \subseteq [0,1]} \left| \frac{T(J, n)}{n} - |J| \right|,$$

where the supremum is taken over all subintervals J of $[0, 1]$ (the length of J is $|J|$) and $T(J, n)$ is the number of points γ_i in J (for $1 \leq i \leq n$). Consider \mathbb{Z}^n with the dot product, and let \mathcal{T} be a sequence of elements in \mathbb{Z}^n . We say that the sequence \mathcal{T} is Δ -homogenously distributed modulo q if for any $a \in \mathbb{Z}^n$, with at least one coordinate coprime to q , the discrepancy of the sequence $\{[a \cdot t]_q/q\}_{t \in \mathcal{T}}$ is at most Δ .

We would like to consider the columns of \mathcal{M} as a set of elements in \mathcal{T} . However, this sequence is not Δ -homogenously distributed modulo q for small Δ : in the case $\mathcal{M} = \mathcal{M}_{\mathbf{h}}$, for example, by construction we have that $\mathbf{h} \cdot \mathbf{f} = \mathbf{g}$ does not distribute homogenously. We therefore define the following weaker notion.

Definition 1 (Weak homogenous distribution). *Let $u \in \mathcal{O}$, $B > 0$, and let \mathcal{T} be a sequence of elements in \mathbb{Z}^n . We say that \mathcal{T} is (B, u) -weakly Δ -homogenously distributed modulo q if for any $a \notin \mathcal{O}u$ such that $\|a\| < B$, with at least one coordinate coprime to q , the discrepancy of the sequence $\{[a \cdot t]_q/q\}_{t \in \mathcal{T}}$ is at most Δ .*

Theorem 1. *For $\mathbf{f}, \mathbf{g} \in R_q$ let $\mathbf{h} = \mathbf{g}/\mathbf{f}$. Suppose that the set of columns of $\mathcal{M}_{\mathbf{h}}$ is (B, \mathbf{f}) -weakly Δ -homogenously distributed modulo q . Let \mathcal{L}' be an $(n + d)$ -dimensional lattice constructed as above. Then with probability at least $1 - (2B - 1)^n((2B + 1)/q + \Delta)^d$ over the columns of $\mathcal{M}_{\mathbf{h}}$, any vector $(x, y) \in \mathcal{L}'$ such that $\|(x, y)\| < B$ satisfies $(x, [xh]_q) = (\alpha f, \alpha g)$ for some $\alpha \in \mathcal{O}$. The result follows for the $(n' + d)$ -dimensional subring lattice \mathcal{L}' where the assumption is taken over $\mathcal{O}_L := \mathcal{O} \cap L$.*

Proof. Suppose $0 \neq (x, y) \in \mathcal{L}'$ such that $x \notin \mathcal{O}f$. Denote by $(xh)_i$ the i -th coefficient of $[xh]_q$, where $0 \leq i < n$. Suppose also that $\|(x, y)\| < B$, which implies that $|x_i| < B$ and $|(xh)_i| < B$. Denote by P the probability that $|(xh)_i| < B$ for some i . By the assumption on \mathbf{h} we get that

$$P < \frac{2B + 1}{q} + \Delta.$$

Recall that the $d \leq n$ rightmost columns in \mathcal{L}' correspond to (a subset of) the coefficients of multiplication by \mathbf{h} . Thus, the probability that $|(xh)_i| < B$ for all the corresponding columns in \mathcal{L}' is $P^d < ((2B + 1)/q + \Delta)^d$, where the probability is taken over the choice of the d columns, which are assumed to be chosen independently.

There are $(2B - 1)^n$ different possibilities for x such that $|x_i| < B$, with $0 \leq i < n$. Hence, the probability over the chosen columns of $\mathcal{M}_{\mathbf{h}}$ in \mathcal{L}' that there exists a lattice point $\|(x, y)\| < B$ with $x \notin \mathcal{O}f$ is at most

$$(2B - 1)^n P^d < (2B - 1)^n \left(\frac{2B + 1}{q} + \Delta \right)^d.$$

□

Theorem 1 considers $\mathcal{M}_{\mathbf{h}}$, and thus also $\mathcal{M}'_{\mathbf{h}}$, as these are the cases of interest in [May99, KF17] respectively. However, it can be generalized to the other matrices \mathcal{M} described above. In general, it is not known that the columns of $\mathcal{M}_{\mathbf{h}}$ are (B, \mathbf{f}) -weakly homogenously distributed for sufficiently small B . Understanding the distribution of \mathbf{h} is extremely important as it underlies the security of NTRU. In general, \mathbf{h} is not uniformly distributed in R_q , as can be seen from a simple information-theoretic argument. Hence, understanding the distribution of \mathbf{f}^{-1} is important in order to understand the distribution of \mathbf{h} . Banks and Shparlinski [BS01] studied how “well spread” the coefficients of \mathbf{f}^{-1} are, that is, whether they “look and behave like random polynomials”. We remark that the desired property on \mathbf{h} may follow from the behaviour of \mathbf{f}^{-1} , but this property is not well formed.

Moreover, it is standard to assume that $\mathbf{h} = \mathbf{g}/\mathbf{f}$ is indistinguishable from random in R_q , see [LATV12]. We remark that this assumption has a strong relation with the weakly homogenous distribution of $\mathcal{M}_{\mathbf{h}}$. Indeed, if the latter is not true, then one can pick a set of small polynomials a and analyze the distribution of $\{[a \cdot t]_q/q\}_{t \in \mathcal{M}_{\mathbf{h}}}$ to distinguish it from from a random \mathbf{h} . Thus, under the indistinguishability assumption, this set of polynomials has to be negligibly small.

We ran experiments with ternary \mathbf{f}, \mathbf{g} and verified that the coefficients of $[\mathbf{f}^{-1}]_q$ equidistribute in \mathbb{Z}_q and that $\mathcal{M}_{\mathbf{h}}$ is (B, \mathbf{f}) -weakly homogenously distributed for sufficiently small B .⁵ For the rest of the paper, we rely on the following assumption.

Assumption 1. *The set of columns $\mathcal{M}_{\mathbf{h}}$ is (B, \mathbf{f}) -weakly $O(q^{-1})$ -homogenously distributed modulo q for $B \ll q$.*

Corollary 1. *Let $\mathbf{f}, \mathbf{g} \in R_q$ and $\mathbf{h} = \mathbf{g}/\mathbf{f}$ be NTRU private and public keys. Let \mathcal{L}' be the projected NTRU lattice of dimension $n + d$ as above, where $d \geq (n \log(2B + 1) + 1)/\log(q/(2B + 1))$. Under Assumption 1, any vector $(x, y) \in \mathcal{L}'$ such that $\|(x, y)\| < B$ satisfies $(x, [xh]_q) = (\alpha f, \alpha g)$ for some $\alpha \in \mathcal{O}$ with probability at least $1/2 + O(q^{-1})$ over the chosen columns of $\mathcal{M}_{\mathbf{h}}$. The result follows for the $(n' + d)$ -dimensional subring lattice \mathcal{L}' where the assumption is taken over $\mathcal{O}_L := \mathcal{O} \cap L$.*

Proof. Using notation from Theorem 1, the probability that there exists a lattice point $\|(x, y)\| < B$

⁵Similarly to [ABD16], though in a different context, in some cases we found that this property holds for $p\mathbf{h}$ for a small prime p (sometimes $2p$). The case $p = 2$ is the case of half integers, where similar a phenomenon was already noted by [ABD16].

with $x \notin \mathcal{O}f$ is at most

$$\frac{(2B+1)^{n+d}}{q^d} + O(q^{-1}) = \frac{2^{(n+d)(\log(2B+1))}}{2^{d \log(q)}} + O(q^{-1}).$$

Setting $d \geq (n \log(2B+1) + 1)/\log(q/(2B+1))$ gives the result. \square

Setting d as required in Corollary 1, observe that the dimension of \mathcal{L}' is $n+d \geq (n \log(q) + 1)/(\log(q/(2B+1)))$. This is similar to the subring lattice in [KF17, Theorem 6]. Thus, Corollary 1 completes the missing details on the validity of the subring attack and, along with [KF17, Theorem 6], gives a complete analysis of the subring attack under Assumption 1. We formalize this result in the following theorem. For the rest of the paper, β denotes the block size in the BKZ [Sch87] algorithm.

Theorem 2 (Adapted from Theorem 6 in [KF17]). *Let $\mathbf{f}, \mathbf{g} \in R_q$ and $\mathbf{h} = \mathbf{g}/\mathbf{f}$ be NTRU private and public keys satisfying Assumption 1. Let $B = \|v\|$ where $v = (\mathbf{f}, \mathbf{g})$ in the full field, $v = (N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{g}))$ in the subfield and $v = (N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{f})\mathbf{h})$ in the subring. Then for*

$$\frac{\beta}{\log \beta} = \frac{2n' \log q}{\log(q/B)^2}$$

one can find a multiple αv for a non-zero $\alpha \in \mathcal{O}$ with probability at least $1/2 + O(q^{-1})$ over the chosen columns of \mathcal{M} .

4.2 Comparing $\|N_{K/L}(\mathbf{g})\|$ and $\|N_{K/L}(\mathbf{f})\mathbf{h}\|$

In this section, we show that when the subring and subfield lattices are fixed to the same dimension via the projection technique, the subring lattice contains smaller vectors. This leads to two main results. For a particular block size β in the BKZ algorithm, a first result is that one can solve a degree- n overstretched NTRU problem in a $2n'$ -dimension lattice with a smaller modulus q using the subring attack. A second result is that for fixed parameters n and q , by further projection, one can solve an overstretched NTRU problem over a smaller lattice dimension using the subring attack.

In Section 3.2 we showed that a small vector in the subfield lattice, most likely the smallest, is of the form $(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{g}))$, while in the subring lattice, there is a small vector of the form $(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{f})\mathbf{h})$. The f part of these vectors is the same. Our interest is therefore in the g part. Moreover, we know that $N_{K/L}(\mathbf{g})$ is an n' -dimensional vector and that $N_{K/L}(\mathbf{f})\mathbf{h}$ is an n -dimensional vector. Our main objective is to show that these two elements have the same Euclidean norm. It then follows that on average the coefficients of $N_{K/L}(\mathbf{g})$ are larger than the coefficients of $N_{K/L}(\mathbf{f})\mathbf{h}$ (since $n > n'$). When we truncate the latter to n' coordinates, its norm becomes smaller than the norm of $N_{K/L}(\mathbf{g})$. Using an assumption on the distribution of the coefficients of $N_{K/L}(\mathbf{g})$, we quantify the difference in size.

More precisely, Theorem 3 shows, with no further assumptions, that when $[K : L] = 2$, the average size of the coefficients of $N_{K/L}(\mathbf{g})$ is expected to be $\sqrt{2}$ times the average size of the coefficients of $N_{K/L}(\mathbf{f})\mathbf{h}$. In Theorem 4, we generalize this result for any subfield such that $[K : L] = r$, and prove

that the ratio of the coefficients is \sqrt{r} under Assumption 2, which we state below. The two main results introduced at the beginning of this section are proven in Corollaries 2 and 3.

We start by stating a simple claim which allows us to obtain the ratio of the coefficients once we assume that the ratio of the Euclidean norms is 1.

Claim 1. *Let f, g be two polynomials with coefficients chosen uniformly at random from the same set. We set $f = (u_1, \dots, u_m)$ and $g = (w_1, \dots, w_n)$. Then, $\mathbb{E}[||f||^2] = \mathbb{E}[||g||^2]$ if and only if the expectation of the square of the coefficients satisfies $\mathbb{E}[u_i^2] = \frac{n}{m} \mathbb{E}[w_i^2]$.*

Proof. We have $\mathbb{E}[||f||^2] = m \mathbb{E}[u_i^2]$ and $\mathbb{E}[||g||^2] = n \mathbb{E}[w_i^2]$. We set these equal and rearrange. \square

In light of this result, our aim is to show that the ratio between $||N_{K/L}(\mathbf{g})||$ and $||N_{K/L}(\mathbf{f})\mathbf{h}||$ tends to 1 as n increases. Then, we can conclude that the subring lattice of dimension $2n'$ is expected to contain shorter vectors than the subfield lattice of the same dimension.

We use random walks to model the coefficients of a product of two polynomials. A first case is for polynomials whose coefficients are drawn independently and uniformly from the set $\{-1, 0, 1\}$. A one-dimensional random walk over \mathbb{Z} starts at 0 and at each step moves either $+1$ or -1 with equal probability. Let a_i , for $i = 1, \dots, n$, denote independent random variables with value either $+1$ or -1 with uniform probability, and let $w_0 = 0$ and $w_n = \sum_{i=1}^n a_i$. The series $\{w_n\}$ defines a random walk over \mathbb{Z} . The expected distance after n steps is on the order of \sqrt{n} . As n increases, the distribution of the series w_n approaches the normal distribution. A second case is for polynomials whose coefficients are drawn from a Gaussian distribution: in a Gaussian random walk, we let a_i follow the Gaussian distribution with standard deviation σ and mean zero. The expected distance after n steps is then on the order of $\sigma\sqrt{n}$. For further background, see [LL10].

Recall that $N_{K/L}(\mathbf{g}) = \prod_{\sigma \in G} \sigma(\mathbf{g})$ and $N_{K/L}(\mathbf{f})\mathbf{h} = \mathbf{g} \prod_{\sigma \in G \setminus \{\text{Id}\}} \sigma(\mathbf{f})$ for $G = \text{Gal}(K/L)$. Moreover, we have $\sigma_1 := \text{Id}$. We now consider the specific case of subfield $L \subseteq K$ of index 2.

Theorem 3. *Let $\mathbf{f}, \mathbf{g} \in R_q$ be two polynomials whose coefficients are drawn independently and uniformly from the set $\{-1, 0, 1\}$, and let $N_{K/L}(\mathbf{g}) = (u_1, \dots, u_n)$ and $\mathbf{g}\sigma_2(\mathbf{f}) = (w_1, \dots, w_n)$. Then $\mathbb{E}[u_i^2 \mid u_i \neq 0] = 8n/9 - 4/9$ and $\mathbb{E}[w_i^2] = 4n/9$. Thus, as n goes to infinity the expected ratio between the non-zero coefficients of $N_{K/L}(\mathbf{g})$ and $\mathbf{g}\sigma_2(\mathbf{f})$ tends to $\sqrt{2}$ in absolute value. In addition, the ratio of the expected squared Euclidean norms $\mathbb{E}[||N_{K/L}(\mathbf{g})||^2] / \mathbb{E}[||\mathbf{g}\sigma_2(\mathbf{f})||^2]$ tends to 1 as n goes to infinity.*

Proof. We start by comparing the size of the coefficients of $N_{K/L}(\mathbf{g})$ to those of $\mathbf{g}\sigma_2(\mathbf{f})$. Let $a_i \in \{-1, 0, 1\}$ and consider the polynomial

$$\mathbf{g} = a_{n-1}x^{n-1} + a_{n-2}x^{n-2} + \dots + a_1x + a_0 \in R_q.$$

Then $\sigma_2(\mathbf{g}) = -a_{n-1}x^{n-1} + a_{n-2}x^{n-2} + \dots - a_1x + a_0$. Observe that each coefficient u_k of $N_{K/L}(\mathbf{g}) = \mathbf{g}\sigma_2(\mathbf{g})$ is a sum of n terms. For k odd, $u_k = 0$, because each of the terms $a_i a_j$ in this sum appears twice with opposite signs. For k even, $u_k = 2 \sum_{i < j, i+j \equiv k \pmod n} \pm a_i a_j + a_{k/2}^2 + a_{(n+k)/2}^2$, since each term $a_i a_j$ with $i \neq j$ appears twice with similar sign. Since $a_i \in \{-1, 0, 1\}$ uniformly and independently

at random, we have that $E[a_i a_j \mid i \neq j] = 0$, $E[a_i^2] = E[a_i^4] = 2/3$, and $E[(a_i a_j)^2] = 4/9$. Thus for k even we have

$$\begin{aligned} E[u_k^2] &= E\left[2 \sum_{i < j, i+j \equiv k \pmod n} \pm a_i a_j + a_{k/2}^2 + a_{(n+k)/2}^2\right] = 4 \sum_{i < j, i+j \equiv k \pmod n} E[(a_i a_j)^2] + E[a_{k/2}^4] + E[a_{(n+k)/2}^4] \\ &= 4(n/2 - 1)(4/9) + 4/3 = 8n/9 - 4/9. \end{aligned}$$

Now, let us repeat a very similar argument for $\mathbf{g}\sigma_2(\mathbf{f})$, where $\sigma_2(\mathbf{f})$ is a polynomial of the form

$$\sigma_2(\mathbf{f}) = b_{n-1}x^{n-1} + b_{n-2}x^{n-2} + \dots + b_1x + b_0,$$

with $b_i \in \{-1, 0, 1\}$ uniformly and independently at random. Again, each coefficient w_k of $\mathbf{g}\sigma_2(\mathbf{f})$ is a sum of n terms. However, unlike the case above, there are no similar terms in this sum, and we have $E[a_i b_j] = 0$ and $E[(a_i b_j)^2] = 4/9$. As above, we compute

$$E[w_k^2] = E\left[\left(\sum_{i+j \equiv k \pmod n} a_i b_j\right)^2\right] = \sum_{i+j \equiv k \pmod n} E[(a_i b_j)^2] = 4n/9.$$

Thus we can compute the expected square of the Euclidean norms

$$\begin{aligned} E[||N_{K/L}(\mathbf{g})||^2] &= E\left[\sum_i u_i^2\right] = \sum_i E[u_i^2] = (n/2)(8n/9 - 4/9) = 4n^2/9 - 2n/9 \\ E[||\mathbf{g}\sigma_2(\mathbf{f})||^2] &= E\left[\sum_i w_i^2\right] = \sum_i E[w_i^2] = n(4n/9) = 4n^2/9. \end{aligned}$$

The result follows. □

We would like to generalize this result to $r > 2$. While the coefficients of $\mathbf{g}\sigma_2(\mathbf{g})$ and $\mathbf{g}\sigma_2(\mathbf{f})$ can be expressed as random walks, and thus follow a Gaussian distribution, they may not be independent. To generalize Theorem 3 we state the following assumption.

Assumption 2. *For $[K : L] > 1$, the coefficients of $N_{K/L}(\mathbf{g})$, $N_{K/L}(\mathbf{f})$ and $N_{K/L}(\mathbf{f})\mathbf{h}$ behave as if they were independently chosen from a Gaussian distribution.*

This assumption seems natural and allows us to prove Theorem 4, a generalisation of Theorem 3 to any $r > 0$. We remark that the results on the attacks given in Corollaries 2 and 3 rely on the result stated in Theorem 4 and not on Assumption 2. We experimentally verified that as n grows, the ratio of the norms tends to 1, as in Theorem 4, and the ratio of the coefficients tends to \sqrt{r} . See Figure 1.

Theorem 4. *Let $\mathbf{f}, \mathbf{g} \in R_q$ and $\mathbf{h} = \mathbf{g}/\mathbf{f}$ be NTRU private and public keys satisfying Assumption 2. For a subfield $L \subseteq K$, let $N_{K/L}(\mathbf{g}) = (u_1, \dots, u_n)$ and $N_{K/L}(\mathbf{f})\mathbf{h} = (w_1, \dots, w_n)$. Then, as n goes to infinity the expected ratio between the non-zero coefficients of $N_{K/L}(\mathbf{g})$ and $N_{K/L}(\mathbf{f})\mathbf{h}$ tends to \sqrt{r} in absolute value. In addition, the ratio of the expected squared Euclidean norms $E[||N_{K/L}(\mathbf{g})||^2]/E[||N_{K/L}(\mathbf{f})\mathbf{h}||^2]$ tends to 1 as n goes to infinity.*

Proof. We give a proof by induction on the index r . The base case is proven in Theorem 3. The general case is a straightforward generalization. Suppose the claim holds for index r , we show that it holds for $[K : L] = 2r$.

First note that for a tower of fields $L \subseteq E \subseteq K$ we have $N_{K/L}(\mathbf{a}) = N_{E/L}(N_{K/E}(\mathbf{a}))$ for every $\mathbf{a} \in K$ (see [LN97, Theorem 2.29]). Consider the case $[K : E] = r, [E : L] = 2$, and denote $\mathbf{G} := N_{K/E}(\mathbf{g})$ and $\mathbf{F} := N_{K/E}(\mathbf{f})$. Then,

$$N_{K/L}(\mathbf{g}) = N_{E/L}(\mathbf{G}) = \mathbf{G}\sigma'_2(\mathbf{G}) \quad \text{and} \quad N_{K/L}(\mathbf{f})\mathbf{h} = N_{E/L}(\mathbf{F})\mathbf{h} = \mathbf{F}\sigma'_2(\mathbf{F})\mathbf{h} = \mathbf{G}'\sigma'_2(\mathbf{F}),$$

where $\sigma'_2 \in \text{Gal}(E/L)$ and $\mathbf{G}' = \mathbf{F}\mathbf{h} = N_{K/E}(\mathbf{f})\mathbf{h}$.

The previous case of the construction, i.e. $[K : E] = r$, shows that each (non-zero) coefficient of \mathbf{F}, \mathbf{G} and \mathbf{G}' follow a Gaussian distribution. Under Assumption 2 the coefficients can be considered to be independent. We can now repeat the process of Theorem 3, which we briefly explain.

While $\mathbf{G}' \in K$, note that $\mathbf{F}, \mathbf{G} \in E$ so they have n/r non-zero coefficients. Thus, similarly to Theorem 3, each non-zero coefficient of $\mathbf{G}\sigma'_2(\mathbf{G})$ is approximately 2 multiplied by a Gaussian random walk with $n/2r$ steps, while each coefficient of $\mathbf{G}'\sigma'_2(\mathbf{F})$ is a random walk with n/r steps. By the induction hypothesis, a coefficient of \mathbf{G} is expected to be larger than the coefficients of \mathbf{G}' by a factor that approaches \sqrt{r} for sufficiently large n . The result on the coefficients follows from evaluating the expected size of the random walks. Then, the claim on the norms follows from Claim 1. \square

As a corollary, we obtain the following result that shows that if we compare subfield and subring attack lattices of the same dimension, then the subring lattice contains smaller vectors.

Corollary 2. *Let $B_{\text{subfield}} = \|(N_{K/L}(\mathbf{f}), N_{K/L}(\mathbf{g}))\|$ where $[K : L] = r = n/n'$. Let $B_{\text{subring}}^{\text{proj}} = \|(N_{K/L}(\mathbf{f}), \overline{N_{K/L}(\mathbf{f})\mathbf{h}})\|$ where the projection keeps $2n'$ coordinates. Under Assumptions 1 and 2, for sufficiently large n we expect to have*

$$\frac{B_{\text{subfield}}^2}{\left(B_{\text{subring}}^{\text{proj}}\right)^2} \approx \frac{2r}{r+1}.$$

Moreover, suppose that B_{subfield} is the norm of the shortest vector in the subfield lattice, and denote by q_{subfield} and q_{subring} the modulus in the subfield and subring attacks, respectively, that the BKZ algorithm can solve NTRU with a fixed block size β . Then

$$\frac{q_{\text{subfield}}}{q_{\text{subring}}} \geq \frac{2r}{r+1}.$$

Proof. To simplify notations, we will write $\text{coeff}(f)$ to denote the coefficients of a polynomial f . We have $B_{\text{subfield}}^2 = (n/r)(\text{coeff}(N_{K/L}(\mathbf{f}))^2 + \text{coeff}(N_{K/L}(\mathbf{g}))^2)$, and

$$\left(B_{\text{subring}}^{\text{proj}}\right)^2 = (n/r)(\text{coeff}(N_{K/L}(\mathbf{f}))^2 + \text{coeff}(N_{K/L}(\mathbf{f})\mathbf{h})^2).$$

We know that $\text{coeff}(N_{K/L}(\mathbf{f}))^2 \approx \text{coeff}(N_{K/L}(\mathbf{g}))^2$ and from Theorem 4, we know that

$$\text{coeff}(N_{K/L}(\mathbf{g}))^2 \approx r \text{coeff}(N_{K/L}(\mathbf{f})\mathbf{h})^2.$$

BKZ is guaranteed to output a vector bounded by $\beta^{2n'/\beta} \lambda_1(\mathcal{L})$. The second result follows from bounding this value by \sqrt{q} in both lattices. \square

It follows from the corollary that taking the same lattice dimension in the subring and subfield attacks, the subring lattice contains vectors of smaller size. Therefore one can solve the NTRU problem using the subring attack with a smaller q . As mentioned in [ABD16, Section 6] we remark that it is not known that $B_{subfield}$ is indeed the norm of the shortest vector (see also [KF17, Theorem 5]). Moreover, our experiments in Table 2 show that the ratio between $q_{subfield}$ and $q_{subring}$ grows with r . A possible explanation is the following: the vector $(N_{K/L}(\mathbf{f}), \overline{N_{K/L}(\mathbf{f})\mathbf{h}})$ is unbalanced, as its f part is shown to be much larger than the g part. Therefore, if there exists an integral multiple of this vector that decreases the size of the f part and increases the size of the g part so that the vector becomes balanced, then the ratio between the feasible qs would increase.

If the systems of equations derived from the lattices given in Corollary 2 are overdetermined, it is possible to project and obtain lower dimensional lattices. Since the g part in the subring lattice is smaller than the subfield lattice, the system in the subring attack is more determined than the system in the subfield attack. That is, the subring lattice will contain solutions of smaller size. Thus, one can discard more equations in the subring attack and achieve a lower dimension. See Table 3 for this comparison. Following Theorem 2, this result gives a better performance for the subring attack as we now show in Corollary 3. Observe that one can get a tighter bound by decreasing the bound B as we project.

Corollary 3. *With the notation from Theorem 2 and Corollary 2, set $B := B_{subfield}$. Under Assumptions 1 and 2, for sufficiently large n , we can find a multiple αv for some non-zero $\alpha \in \mathcal{O}$ such that using BKZ with block size β on the lattice \mathcal{L}' of dimension $n' + d$,*

1. *if \mathcal{L}' is the subfield lattice, then*

$$n' + d \geq \frac{n' \log(q) + 1}{\log(q/(2B + 1))} \quad \text{and} \quad \frac{\beta}{\log \beta} = \frac{2n' \log q}{(\log(q/B))^2},$$

2. *if \mathcal{L}' is the subring lattice, then*

$$n' + d \geq \frac{n' \log(q) + 1}{\log(q\sqrt{r}/(\sqrt{2r} + 2B + \sqrt{r}))} \quad \text{and} \quad \frac{\beta}{\log \beta} = \frac{2n' \log q}{(\log(\sqrt{2r}q/\sqrt{r} + 1B))^2}.$$

4.2.1 Experimental results

We implemented the full field, subfield, and subring attacks in Sage and experimentally compared them using Sage's default LLL implementation. We defined a success in our experiments as recovering a vector v satisfying $\|v\| < q^{3/4}$. As already noted by [KF17], we either get vectors which are roughly of size q , or vectors of size \sqrt{q} that are short integral multiples of the private key.

We fix the parameters $(n, \text{dimension}, r)$ and compare the smallest modulus q that succeeded for each of the full field, subring, and subfield attacks using LLL. In all cases, the subring attack succeeded for a smaller modulus than the subfield attack, and outperformed our analytic bounds in the larger experiments, see Table 2. Lattice reduction for the 512-dimensional lattice for the

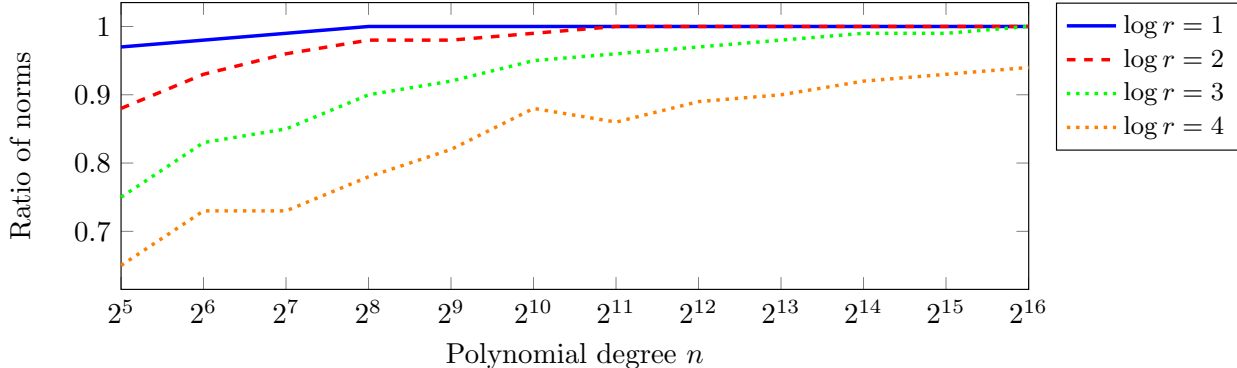


Figure 1: **Experimentally verifying that the ratio of the norms converges to 1.** We experimentally computed the ratio $\|N_{K/L}(\mathbf{g})\|/\|N_{K/L}(\mathbf{f})\mathbf{h}\|$ for n increasing and for various subfields of index r . The ratio converges experimentally to 1 as $n \rightarrow \infty$ for each index r we experimented with. The speed of the convergence is proportional to r .

subring attack for $\log n = 12$ took 276 hours of wall clock time and the 1024-dimensional lattice reduction for the full field attack with $\log n = 9$ took 182 hours of wall clock time on a single core of our experimental machine.

We then compared the subring and subfield attacks by fixing (n, q, r) and comparing the smallest dimension that succeeded using LLL. For some of the lattices we used the projection technique by deleting the right-most columns until we reached the desired dimension. For $n \in [2^6, 2^{11}]$ and $r \in [2, 2^4]$, the subring attack succeeded with a smaller dimension than the subfield attack. In Table 3, we give detailed experimental results. We also report timings for lattice reduction done with LLL. The longest computations in the table were the last few rows; lattice reduction took 4.2 hours of wall clock time for the 254-dimensional lattice on a single core of an Intel Xeon E5-2699 v3 running at 2.30GHz, with 128 GB of RAM. We note that in most cases, for a given dimension, the lattice reduction for the subring attack seems to be slightly more efficient than for the subfield attack.

5 Open questions

In this paper we give a comparison between the subfield and subring attacks on overstretched NTRU. We argue that the correct method for comparison should be the resulting lattice dimension, contrary to the comparisons made in previous work. Our analysis shows that for a fixed dimension, the subring lattice using the projection method contains shorter vectors, and this can be used to solve the NTRU problem with the subring attack on a lattice of smaller dimension than the subfield attack.

As pointed out in previous work, the desired vectors in the lattices are not known to be the smallest, and as the experiments show in some cases they are not. In order to understand the actual difference between the attacks, we need a better understanding of the smallest vectors. We mentioned above that the desired vector in the subring lattice is unbalanced, and therefore

it is theoretically possible that a balanced multiple of it exists. Proving this claim would give insights about the results of our experiments. It also of great interest to explain our experimental observations of cases where the full field attack only succeeds after projection and is outperformed by the other attacks.

The subfield attack in [CJL16] uses the trace function and the g part of the desired lattice vector in this attack appears to be somewhat larger than in the subfield attack that uses the norm. It is not immediately clear what is the ratio between these g parts. The g part in the trace attack is given by $N_{K/L}(\mathbf{f})Tr_{K/L}(\mathbf{h}) = \sum_i \sigma_i(\mathbf{g}) \prod_{j \neq i} \sigma_j(\mathbf{f})$. It seems that our analysis can be used to analyze this case, which we leave to future work.

An interesting question revolves around the existence of shifts in the NTRU lattice. On the one hand, their existence is fundamental for Kirchner and Fouque’s analysis of the attack in the full field [KF17, Theorem 9]. On the other hand, the goal of the zero forcing method is to eliminate them. We ran some experiments to check whether eliminating the shifts is a better approach than removing the corresponding rows, and observe that both methods lead to the same result. The only benefit of removing some rows would be the slight improvement in the runtime due to a smaller lattice dimension. We note that some other related interesting open questions are given in [MS01].

Acknowledgements We thank Shi Bai for some helpful discussions on the topic. This material is based upon work supported by the National Science Foundation under Grants No. CNS-1513671 and CNS-1651344 and by the ONR SynCrypt project. We are grateful to Cisco for donating the Cisco UCS servers that we used for our experiments.

References

- [ABD16] Martin Albrecht, Shi Bai, and Léo Ducas. A Subfield Lattice Attack on Overstretched NTRU Assumptions. In Matthew Robshaw and Jonathan Katz, editors, *Advances in Cryptology – CRYPTO 2016*, pages 153–178, Berlin, Heidelberg, 2016. Springer Berlin Heidelberg.
- [BLLN13] Joppe W. Bos, Kristin Lauter, Jake Loftus, and Michael Naehrig. Improved security for a ring-based fully homomorphic encryption scheme. In Martijn Stam, editor, *Cryptography and Coding*, pages 45–64, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [BS01] William D. Banks and Igor E. Shparlinski. Distribution of Inverses in Polynomial Rings. *Indagationes Mathematicae*, 12(3), 2001.
- [CJL16] Jung Hee Cheon, Jinhyuck Jeong, and Changmin Lee. An Algorithm for NTRU Problems and Cryptanalysis of the GGH Multilinear Map without a Low-Level Encoding of Zero. *LMS Journal of Computation and Mathematics*, 19(A):255–266, 2016.

- [CS97] Don Coppersmith and Adi Shamir. Lattice Attacks on NTRU. In Walter Fumy, editor, *Advances in Cryptology — EUROCRYPT '97*, pages 52–61, Berlin, Heidelberg, 1997. Springer Berlin Heidelberg.
- [DT97] M. Drmota and R. Tichy. *Discrepancies and Applications*. Springer-Verlag, Berlin, 1997.
- [DYT17] Dung Hoang Duong, Masaya Yasuda, and Tsuyoshi Takagi. Choosing Parameters for the Subfield Lattice Attack Against Overstretched NTRU. In Phong Q. Nguyen and Jianying Zhou, editors, *Information Security*, pages 79–91, Cham, 2017. Springer International Publishing.
- [GGH13] Sanjam Garg, Craig Gentry, and Shai Halevi. Candidate Multilinear Maps from Ideal Lattices. In Thomas Johansson and Phong Q. Nguyen, editors, *Advances in Cryptology — EUROCRYPT 2013*, pages 1–17, Berlin, Heidelberg, 2013. Springer Berlin Heidelberg.
- [GS02] Craig Gentry and Mike Szydlo. Cryptanalysis of the revised ntru signature scheme. In Lars R. Knudsen, editor, *Advances in Cryptology — EUROCRYPT 2002*, pages 299–320, Berlin, Heidelberg, 2002. Springer Berlin Heidelberg.
- [HHGP⁺03] Jeffrey Hoffstein, Nick Howgrave-Graham, Jill Pipher, Joseph H. Silverman, and William Whyte. NTRUSign: Digital Signatures Using the NTRU Lattice. In Marc Joye, editor, *Topics in Cryptology — CT-RSA 2003*, pages 122–140, Berlin, Heidelberg, 2003. Springer Berlin Heidelberg.
- [HPS96] Jeffrey Hoffstein, Jill Pipher, and Joseph H. Silverman. NTRU: A New High Speed Public Key Cryptosystem, 1996.
- [HPS98] Jeffrey Hoffstein, Jill Pipher, and Joseph H. Silverman. NTRU: A Ring-based Public Key Cryptosystem. In Joe P. Buhler, editor, *Algorithmic Number Theory*, pages 267–288, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg.
- [KF17] Paul Kirchner and Pierre-Alain Fouque. Revisiting Lattice Attacks on Overstretched NTRU Parameters. In Jean-Sébastien Coron and Jesper Buus Nielsen, editors, *Advances in Cryptology — EUROCRYPT 2017*, pages 3–26, Cham, 2017. Springer International Publishing.
- [KN74] R. Kuipers and H. Niederreiter. *Uniform Distribution of Sequences*. Wiley-Interscience, New York, NY, USA, 1974.
- [LATV12] Adriana López-Alt, Eran Tromer, and Vinod Vaikuntanathan. On-the-fly Multiparty Computation on the Cloud via Multikey Fully Homomorphic Encryption. In *Proceedings of the Forty-fourth Annual ACM Symposium on Theory of Computing, STOC '12*, pages 1219–1234, New York, NY, USA, 2012. ACM.

- [LL10] Gregory F. Lawler and Vlada Limic. *Random Walk: A Modern Introduction*. Cambridge University Press, New York, NY, USA, 2010.
- [LN97] R. Lidl and H. Niederreiter. *Finite Fields*. Cambridge University Press, 1997.
- [May99] Alexander May. Cryptanalysis of NTRU, 1999. Preprint.
- [MS01] Alexander May and Joseph H. Silverman. Dimension Reduction Methods for Convolution Modular Lattices. In Joseph H. Silverman, editor, *Cryptography and Lattices*, pages 110–125, Berlin, Heidelberg, 2001. Springer Berlin Heidelberg.
- [PT08] Gábor Pataki and Mustafa Tural. On Sublattice Determinants in Reduced Bases, 2008.
- [Sam70] P. Samuel. *Algebraic Theory of Numbers*. Hermann, 1970.
- [Sch87] C. P. Schnorr. A Hierarchy of Polynomial Time Lattice Basis Reduction Algorithms. *Theor. Comput. Sci.*, 53(2-3):201–224, August 1987.
- [SS11] Damien Stehlé and Ron Steinfeld. Making NTRU as Secure as Worst-Case Problems over Ideal Lattices. In Kenneth G. Paterson, editor, *Advances in Cryptology – EUROCRYPT 2011*, pages 27–47, Berlin, Heidelberg, 2011. Springer Berlin Heidelberg.
- [Ste14] Ron Steinfeld. *NTRU Cryptosystem: Recent Developments and Emerging Mathematical Problems in Finite Polynomial Rings*, pages 179 – 211. Walter de Gruyter, 2014.

Parameters					Attacks		
$\log n$	Dimension (Full/subfield)	$\log r$	$\log q$	Full field	Subring	Subfield	
6	128/64	1	13	Yes	Yes	Yes	
6	128/64	1	11	Yes	Yes	No	
6	128/64	1	9	Yes	No	No	
6	128/64	1	8	No	No	No	
7	256/128	1	16	Yes	Yes	Yes	
7	256/128	1	15	Yes	Yes	No	
7	256/128	1	14	Yes	No	No	
7	256/128	1	13	No	No	No	
8	(384;512)/256	1	22	Yes;No	Yes	Yes	
8	(384;512)/256	1	21	No	Yes	No	
9	(768;1024)/512	1	34	Yes;No	Yes	Yes	
9	768/512	1	32	No	Yes	Yes	
9	768/512	1	31	No	No	No	
9	1024/256	2	40	Yes	Yes	Yes	
9	1024/256	2	38	Yes	Yes	No	
9	1024/256	2	37	Yes	No	No	
9	1024/256	2	36	Yes	No	No	
10	2048/512	2	52	-	Yes	Yes	
10	2048/512	2	50	-	Yes	No	
10	2048/512	2	49	-	No	No	
11	4096/512	3	95	-	Yes	Yes	
11	4096/512	3	92	-	Yes	No	
11	4096/512	3	91	-	No	No	
11	4096/256	4	165	-	Yes	Yes	
11	4096/256	4	162	-	Yes	No	
11	4096/256	4	161	-	No	No	
12	4096/512	4	189	-	Yes	Yes	
12	4096/512	4	185	-	Yes	No	
12	4096/512	4	184	-	No	No	

Table 2: **Experimentally determining the minimal q for each attack.** For a fixed dimension, we compare the subfield and the subring attacks on NTRU with the *same* modulus q and note whether the attack succeeded. In some cases the full field attack only succeeded when we projected to a smaller lattice. We represent this for example with the notation (384;512), which means that we ran the full field attack on a 384-dimensional lattice and on a 512-dimensional lattice.

Parameters				Attacks		LLL reduction time (s)	
$\log n$	$\log q$	$\log r$	Dimension	Subfield	Subring	Subfield	Subring
6	13	1	51	Yes	Yes	0.464	0.483
6	13	1	50	No	Yes	0.409	0.409
6	13	1	49	No	No	0.382	0.373
7	16	1	119	Yes	Yes	18.9	17.5
7	16	1	111	No	Yes	14.0	14.1
7	16	1	110	No	No	13.1	12.8
8	22	1	233	Yes	Yes	715.4	577.6
8	22	1	223	No	Yes	454.5	452.3
8	22	1	222	No	No	456.5	424.3
9	70	3	122	Yes	Yes	132.6	121.1
9	70	3	117	No	Yes	116.2	117.0
9	70	3	116	No	No	108.7	108.9
10	150	4	124	Yes	Yes	344.8	325.3
10	150	4	120	No	Yes	298.6	304.1
10	150	4	119	No	No	290.2	286.5
11	165	4	254	Yes	Yes	15333.0	12423.7
11	165	4	249	No	Yes	12708.7	12072.0
11	165	4	248	No	No	12086.3	11722.8

Table 3: **Experimentally determining the minimal dimension for each attack.** For a fixed modulus q , we compare the subfield and the subring attacks applied to a lattice of the *same* dimension and note whether the attack succeeded. Bold-faced **Yes** indicates the lowest dimension we reached for a given set of parameters. Experimentally, the subring attack succeeds with smaller dimension lattices than the subfield attack.