

Post-quantum key exchange – a new hope*

Erdem Alkim

Department of Mathematics, Ege University, Turkey

Léo Ducas

Centrum voor Wiskunde en Informatica, Amsterdam, The Netherlands

Thomas Pöppelmann

Infineon Technologies AG, Munich, Germany

Peter Schwabe

Digital Security Group, Radboud University, The Netherlands

Abstract

At IEEE Security & Privacy 2015, Bos, Costello, Naehrig, and Stebila proposed an instantiation of Peikert’s ring-learning-with-errors-based (Ring-LWE) key-exchange protocol (PQCrypto 2014), together with an implementation integrated into OpenSSL, with the affirmed goal of providing post-quantum security for TLS. In this work we revisit their instantiation and stand-alone implementation. Specifically, we propose new parameters and a better suited error distribution, analyze the scheme’s hardness against attacks by quantum computers in a conservative way, introduce a new and more efficient error-reconciliation mechanism, and propose a defense against backdoors and all-for-the-price-of-one attacks. By these measures and for the same lattice dimension, we more than double the security parameter, halve the communication overhead, and speed up computation by more than a factor of 8 in a portable C implementation and by more than a factor of 27 in an optimized implementation targeting current Intel CPUs. These speedups are achieved with comprehensive protection against timing attacks.

1 Introduction

The last decade in cryptography has seen the birth of numerous constructions of cryptosystems based on lattice problems, achieving functionalities that were previously unreachable (e.g., fully homomorphic cryptogra-

phy [42]). But even for the simplest tasks in asymmetric cryptography, namely public-key encryption, signatures, and key exchange, lattice-based cryptography offers an important feature: resistance to all known quantum algorithms. In those times of *quantum nervousness* [76, 77], the time has come for the community to deliver and optimize concrete schemes, and to get involved in the standardization of a lattice-based cipher-suite via an open process.

For encryption and signatures, several competitive schemes have been proposed; examples are NTRU encryption [54, 87], Ring-LWE encryption [70] as well as the signature schemes BLISS [35], PASS [52] or the proposal by Bai and Galbraith presented in [8]. To complete the lattice-based cipher-suite, Bos, Costello, Naehrig, and Stebila [22] recently proposed a concrete instantiation of Peikert’s improved version of the key-exchange scheme originally proposed by Ding, Xie, and Lin [34, 80]. Bos et al. proved its practicality by integrating their implementation as additional cipher-suite into the transport layer security (TLS) protocol in OpenSSL. In the following we will refer to this proposal as BCNS.

Unfortunately, the performance of BCNS seemed rather disappointing. We identify two main sources for this inefficiency. First the analysis of the failure probability was far from tight, resulting in a very large modulus $q \approx 2^{32}$. As a side effect, the security is also significantly lower than what one could achieve with Ring-LWE for a ring of rank $n = 1024$. Second the Gaussian sampler, used to generate the secret parameters, is fairly inefficient and hard to protect against timing attacks. This second source of inefficiency stems from the fundamental misconception that high-quality Gaussian noise is crucial for encryption based on LWE¹, which has also made various other implementations [32, 82] slower and more complex than they would have to be.

*This work was initiated while Thomas Pöppelmann was a Ph.D. student at Ruhr-University Bochum with support from the European Union H2020 SAFEcrypto project (grant no. 644729). This work has furthermore been supported by TÜBITAK under 2214-A Doctoral Research Program Grant, by the European Commission through the ICT program under contract ICT-645622 (PQCRYPTO), and by the Netherlands Organisation for Scientific Research (NWO) through Veni 2013 project 13114 and through a Free Competition Grant. Permanent ID of this document: 0462d84a3d34b12b75e8f5e4ca032869. Date: 2016-08-03.

¹This is very different for lattice-based signatures or trapdoors, where distributions need to be meticulously crafted to prevent any leak of information on a secret basis.

1.1 Contributions

In this work, we propose solutions to the performance and security issues of the aforementioned BCNS proposal [22]. Our improvements are possible through a combination of multiple contributions:

- Our first contribution is an improved analysis of the failure probability of the protocol. To push the scheme even further, inspired by analog error-correcting codes, we make use of the lattice D_4 to allow error reconciliation beyond the original bounds of [80]. This drastically decreases the modulus to $q = 12289 < 2^{14}$, which improves both efficiency and security.
- Our second contribution is a more detailed security analysis against quantum attacks. We provide a lower bound on all known (or even presupposed) quantum algorithms solving the shortest-vector problem (SVP), and deduce the potential performance of a quantum BKZ algorithm. According to this analysis, our improved proposal provides 128 bits of post-quantum security with a comfortable margin.
- We furthermore propose to replace the almost-perfect discrete Gaussian distribution by something relatively close, but much easier to sample, and prove that this can only affect the security marginally.
- We replace the fixed parameter \mathbf{a} of the original scheme by a freshly chosen random one in each key exchange. This incurs an acceptable overhead but prevents backdoors embedded in the choice of this parameter and all-for-the-price-of-one attacks.
- We specify an encoding of polynomials in the number-theoretic transform (NTT) domain which allows us to eliminate some of the NTT transformations inside the protocol computation.
- To demonstrate the applicability and performance of our design we provide a portable reference implementation written in C and a highly optimized vectorized implementation that targets recent Intel CPUs and is compatible with recent AMD CPUs. We describe an efficient approach to lazy reduction inside the NTT, which is based on a combination of Montgomery reductions and short Barrett reductions.

Availability of software. We place all software described in this paper into the public domain and make it available online at <https://cryptojedi.org/crypto/#newhope> and <https://github.com/tpoeppelmann/newhope>.

Acknowledgments. We are thankful to Mike Hamburg and to Paul Crowley for pointing out mistakes in a previous version of this paper, and we are thankful to Isis Lovecruft for thoroughly proofreading the paper and for suggesting the name JARJAR for the low-security variant of our proposal.

2 Lattice-based key exchange

Let \mathbb{Z} be the ring of rational integers. We define for an $x \in \mathbb{R}$ the rounding function $\lfloor x \rfloor = \lfloor x + \frac{1}{2} \rfloor \in \mathbb{Z}$. Let \mathbb{Z}_q , for an integer $q \geq 1$, denote the quotient ring $\mathbb{Z}/q\mathbb{Z}$. We define $\mathcal{R} = \mathbb{Z}[X]/(X^n + 1)$ as the ring of integer polynomials modulo $X^n + 1$. By $\mathcal{R}_q = \mathbb{Z}_q[X]/(X^n + 1)$ we mean the ring of integer polynomials modulo $X^n + 1$ where each coefficient is reduced modulo q . In case χ is a probability distribution over \mathcal{R} , then $x \stackrel{\$}{\leftarrow} \chi$ means the sampling of $x \in \mathcal{R}$ according to χ . When we write $\mathbf{a} \stackrel{\$}{\leftarrow} \mathcal{R}_q$ this means that all coefficients of \mathbf{a} are chosen uniformly at random from \mathbb{Z}_q . For a probabilistic algorithm \mathcal{A} we denote by $y \stackrel{\$}{\leftarrow} \mathcal{A}$ that the output of \mathcal{A} is assigned to y and that \mathcal{A} is running with randomly chosen coins. We recall the discrete Gaussian distribution $D_{\mathbb{Z}, \sigma}$ which is parametrized by the Gaussian parameter $\sigma \in \mathbb{R}$ and defined by assigning a weight proportional to $\exp(-\frac{x^2}{2\sigma^2})$ to all integers x .

2.1 The scheme of Peikert

In this section we briefly revisit the passively secure key-encapsulation mechanism (KEM) that was proposed by Peikert [80] and instantiated in [22] (BCNS). Peikert’s KEM scheme is defined by the algorithms (Setup, Gen, Encaps, Decaps) and after a successful protocol run both parties share an ephemeral secret key that can be used to protect further communication (see Protocol 1).

The KEM scheme by Peikert closely resembles a previously introduced Ring-LWE encryption scheme [69] but due to a new error-reconciliation mechanism, one \mathcal{R}_q component of the ciphertext can be replaced by a more compact element in \mathcal{R}_2 . This efficiency gain is possible due to the observation that it is not necessary to transmit an explicitly chosen key to establish a secure ephemeral session key. In Peikert’s scheme, the reconciliation just allows both parties to derive the session key from an approximately agreed pseudorandom ring element. For Alice, this ring element is $\mathbf{u} = \mathbf{a}\mathbf{s}' + \mathbf{e}'\mathbf{s}$ and for Bob it is $\mathbf{v} = \mathbf{b}\mathbf{s}' + \mathbf{e}'' = \mathbf{a}\mathbf{s}' + \mathbf{e}'\mathbf{s} + \mathbf{e}''$. For a full explanation of the reconciliation we refer to the original paper [80] but briefly recall the cross-rounding function $\langle \cdot \rangle_2$ defined as $\langle v \rangle_2 := \lfloor \frac{4}{q} \cdot v \rfloor \bmod 2$ and the randomized function $\text{dbl}(v) := 2v - \bar{e}$ for some random \bar{e} where $\bar{e} = 0$ with probability $\frac{1}{2}$, $\bar{e} = 1$ with probability $\frac{1}{4}$, and

$\bar{e} = -1$ with probability $\frac{1}{4}$. Let $I_0 = \{0, 1, \dots, \lfloor \frac{q}{2} \rfloor - 1\}$, $I_1 = \{-\lfloor \frac{q}{2} \rfloor, \dots, -1\}$, and $E = [-\frac{q}{4}, \frac{q}{4})$ then the reconciliation function $\text{rec}(w, b)$ is defined as

$$\text{rec}(w, b) = \begin{cases} 0, & \text{if } w \in I_b + E \pmod{q} \\ 1, & \text{otherwise.} \end{cases}$$

If these functions are applied to polynomials this means they are applied to each of the coefficients separately.

Parameters: q, n, χ	
KEM.Setup():	
$\mathbf{a} \xleftarrow{\$} \mathcal{R}_q$	
Alice (server)	Bob (client)
KEM.Gen(\mathbf{a}):	KEM.Encaps(\mathbf{a}, \mathbf{b}):
$\mathbf{s}, \mathbf{e} \xleftarrow{\$} \chi$	$\mathbf{s}', \mathbf{e}', \mathbf{e}'' \xleftarrow{\$} \chi$
$\mathbf{b} \leftarrow \mathbf{a}\mathbf{s} + \mathbf{e}$	$\xrightarrow{\mathbf{b}}$
	$\mathbf{u} \leftarrow \mathbf{a}\mathbf{s}' + \mathbf{e}'$
	$\mathbf{v} \leftarrow \mathbf{b}\mathbf{s}' + \mathbf{e}''$
	$\bar{\mathbf{v}} \xleftarrow{\$} \text{dbl}(\mathbf{v})$
KEM.Decaps($\mathbf{s}, (\mathbf{u}, \mathbf{v}')$):	$\mathbf{v}' = \langle \bar{\mathbf{v}} \rangle_2$
$\mu \leftarrow \text{rec}(2\mathbf{u}\mathbf{s}, \mathbf{v}')$	$\mu \leftarrow \lfloor \bar{\mathbf{v}} \rfloor_2$

Protocol 1: Peikert’s KEM mechanism.

2.2 The BCNS proposal

In a work by Bos, Costello, Naehrig, and Stebila [22] (BCNS), Peikert’s KEM [80] was phrased as a key-exchange protocol (see again Protocol 1), instantiated for a concrete parameter set, and integrated into OpenSSL (see Section 8 for a performance comparison). Selection of parameters was necessary as Peikert’s original work does not contain concrete parameters and the security as well as error estimation are based on asymptotics. The authors of [22] chose a dimension $n = 1024$, a modulus $q = 2^{32} - 1$, $\chi = D_{\mathbb{Z}_n, \sigma}$ and the Gaussian parameter $\sigma = 8/\sqrt{2\pi} \approx 3.192$. It is claimed that these parameters provide a classical security level of at least 128 bits considering the distinguishing attack [65] with distinguishing advantage less than 2^{-128} and $2^{81.9}$ bits of security against an optimistic instantiation of a quantum adversary. The probability of a wrong key being established is less than $2^{-2^{17}} = 2^{-131072}$. The message \mathbf{b} sent by Alice is a ring element and thus requires at least $\log_2(q)n = 32$ kbits while Bob’s response (\mathbf{u}, \mathbf{r}) is a ring element R_q and an element from R_2 and thus requires at least 33 kbits. As the polynomial $\mathbf{a} \in \mathcal{R}_q$ is shared between all parties this ring element has to be stored or generated on-the-fly. For timings of their implementation we refer to Table 2. We would also like to note that besides its aim for securing

classical TLS, the BCNS protocol has already been proposed as a building block for Tor [88] on top of existing elliptic-curve infrastructure [45].

2.3 Our proposal: NEWHOPE

In this section we detail our proposal and modifications of Peikert’s protocol². For the same reasons as described in [22] we opt for an unauthenticated key-exchange protocol; the protection of stored transcripts against future decryption using quantum computers is much more urgent than post-quantum authentication. Authenticity will most likely be achievable in the foreseeable future using proven pre-quantum signatures and attacks on the signature will not compromise previous communication. Additionally, by not designing or instantiating a lattice-based authenticated key-exchange protocol (see [37, 90]) we reduce the complexity of the key-exchange protocol and simplify the choice of parameters. We actually see it as an advantage to decouple key exchange and authentication as it allows a protocol designer to choose the optimal algorithm for both tasks (e.g., an ideal-lattice-based key exchange and a hash-based signature like [18] for authentication). Moreover, this way the design, security level, and parameters of the key-exchange scheme are not constrained by requirements introduced by the authentication part.

Parameter choices. A high-level description of our proposal is given in Protocol 2 and as in [22, 80] all polynomials except for $\mathbf{r} \in \mathcal{R}_4$ are defined in the ring $\mathcal{R}_q = \mathbb{Z}_q[X]/(X^n + 1)$ with $n = 1024$ and $q = 12289$. We decided to keep the dimension $n = 1024$ as in [22] to be able to achieve appropriate long-term security. As polynomial arithmetic is fast and also scales better (doubling n roughly doubles the time required for a polynomial multiplication), our choice of n appears to be acceptable from a performance point of view. We chose the modulus $q = 12289$ as it is the smallest prime for which it holds that $q \equiv 1 \pmod{2n}$ so that the number-theoretic transform (NTT) can be realized efficiently and that we can transfer polynomials in NTT encoding (see Section 7).

As the security level grows with the noise-to-modulus ratio, it makes sense to choose the modulus as small as possible, improving compactness and efficiency together with security. The choice is also appealing as the prime is already used by some implementations of Ring-LWE encryption [32, 66, 85] and BLISS signatures [35, 81]; thus sharing of some code (or hardware modules) between our proposal and an implementation of BLISS would be possible.

²For the TLS use-case and for compatibility with BCNS [22] the key exchange is initiated by the server. However, in different scenarios the roles of the server and client can be exchanged.

Noise distribution and reconciliation. Notably, we also change the distribution of the LWE secret and error and replace discrete Gaussians by the centered binomial distribution ψ_k of parameter $k = 16$ (see Section 4). The reason is that it turned out to be challenging to implement a discrete Gaussian sampler efficiently *and* protected against timing attacks (see [22] and Section 5). On the other hand, sampling from the centered binomial distribution is easy and does not require high-precision computations or large tables as one may sample from ψ_k by computing $\sum_{i=0}^k b_i - b'_i$, where the $b_i, b'_i \in \{0, 1\}$ are uniform independent bits. The distribution ψ_k is centered (its mean is 0), has variance $k/2$ and for $k = 16$ this gives a standard deviation of $\zeta = \sqrt{16/2}$. Contrary to [22, 80] we hash the output of the reconciliation mechanism, which makes a distinguishing attack irrelevant and allows us to argue security for the modified error distribution.

Moreover, we generalize Peikert’s reconciliation mechanism using an analog error-correction approach (see Section 5). The design rationale is that we only want to transmit a 256-bit key but have $n = 1024$ coefficients to encode data into. Thus we encode one key bit into four coefficients; by doing so we achieve increased error resilience which in turn allows us to use larger noise for better security.

Short-term public parameters. NEWHOPE does not rely on a globally chosen public parameter \mathbf{a} as the efficiency increase in doing so is not worth the measures that have to be taken to allow trusted generation of this value and the defense against backdoors [15]. Moreover, this approach avoids the rather uncomfortable situation that all connections rely on a single instance of a lattice problem (see Section 3) in the flavor of the “Logjam” DLP attack [1].

No key caching. For ephemeral Diffie-Hellman key-exchange in TLS it is common for servers to cache a key pair for a short time to increase performance. For example, according to [26], Microsoft’s SChannel library caches ephemeral keys for 2 hours. We remark that for the lattice-based key exchange described in [80], for the key exchange described in [22], and also for the key exchange described in this paper, such short-term caching would be disastrous for security. Indeed, it is crucial that both parties use fresh secrets for each instantiation (thus the performance of the noise sampling is crucial). As short-term key caching typically happens on higher layers of TLS libraries than the key-exchange implementation itself, we stress that particular care needs to be taken to eliminate such caching when switching from ephemeral (elliptic-curve) Diffie-Hellman key exchange to post-quantum lattice-based key exchange. This issue is discussed in more detail in [36].

One *could* enable key caching with a transformation from the CPA-secure key exchange to a CCA-secure key exchange as outlined by Peikert in [80, Section 5]. Note that such a transform would furthermore require changes to the noise distribution to obtain a failure probability that is negligible in the cryptographic sense.

3 Preventing backdoors and all-for-the-price-of-one attacks

One serious concern about the original design [22] is the presence of the polynomial \mathbf{a} as a fixed system parameter. As described in Protocol 2, our proposal includes pseudorandom generation of this parameter for every key exchange. In the following we discuss the reasons for this decision.

Backdoor. In the worst scenario, the fixed parameter \mathbf{a} could be backdoored. For example, inspired by NTRU trapdoors [54, 87], a dishonest authority may choose mildly small \mathbf{f}, \mathbf{g} such that $\mathbf{f} = \mathbf{g} = 1 \pmod p$ for some prime $p \geq 4 \cdot 16 + 1$ and set $\mathbf{a} = \mathbf{g}\mathbf{f}^{-1} \pmod q$. Then, given $(\mathbf{a}, \mathbf{b} = \mathbf{a}\mathbf{s} + \mathbf{e})$, the attacker can compute $\mathbf{b}\mathbf{f} = \mathbf{a}\mathbf{f}\mathbf{s} + \mathbf{f}\mathbf{e} = \mathbf{g}\mathbf{s} + \mathbf{f}\mathbf{e} \pmod q$, and, because $\mathbf{g}, \mathbf{s}, \mathbf{f}, \mathbf{e}$ are small enough, compute $\mathbf{g}\mathbf{s} + \mathbf{f}\mathbf{e}$ in \mathbb{Z} . From this he can compute $\mathbf{t} = \mathbf{s} + \mathbf{e} \pmod p$ and, because the coefficients of \mathbf{s} and \mathbf{e} are smaller than 16, their sums are in $[-2 \cdot 16, 2 \cdot 16]$: knowing them modulo $p \geq 4 \cdot 16 + 1$ is knowing them in \mathbb{Z} . It now only remains to compute $(\mathbf{b} - \mathbf{t}) \cdot (\mathbf{a} - 1)^{-1} = (\mathbf{a}\mathbf{s} - \mathbf{s}) \cdot (\mathbf{a} - 1)^{-1} = \mathbf{s} \pmod q$ to recover the secret \mathbf{s} .

One countermeasure against such backdoors is the “nothing-up-my-sleeve” process, which would, for example, choose \mathbf{a} as the output of a hash function on a common universal string like the digits of π . Yet, even this process may be partially abused [15], and when not strictly required it seems preferable to avoid it.

All-for-the-price-of-one attacks. Even if this common parameter has been honestly generated, it is still rather uncomfortable to have the security of all connections rely on a single instance of a lattice problem. The scenario is an entity that discovers an unforeseen cryptanalytic algorithm, making the required lattice reduction still very costly, but say, not impossible in a year of computation, given its outstanding computational power. By finding *once* a good enough basis of the lattice $\Lambda = \{(a, 1)x + (q, 0)y \mid x, y \in \mathcal{R}\}$, this entity could then compromise *all* communications, using for example Babai’s decoding algorithm [7].

This idea of massive precomputation that is only dependent on a fixed parameter \mathbf{a} and then afterwards can be used to break all key exchanges is similar in flavor to the 512-bit “Logjam” DLP attack [1]. This attack was only possible in the required time limit because

Parameters: $q = 12289 < 2^{14}$, $n = 1024$	
Error distribution: ψ_{16}	
Alice (server)	Bob (client)
$seed \xleftarrow{\$} \{0, 1\}^{256}$	
$\mathbf{a} \leftarrow \text{Parse}(\text{SHAKE-128}(seed))$	
$\mathbf{s}, \mathbf{e} \xleftarrow{\$} \psi_{16}^n$	$\mathbf{s}', \mathbf{e}', \mathbf{e}'' \xleftarrow{\$} \psi_{16}^n$
$\mathbf{b} \leftarrow \mathbf{a}\mathbf{s} + \mathbf{e}$	$\mathbf{a} \leftarrow \text{Parse}(\text{SHAKE-128}(seed))$
	$\mathbf{u} \leftarrow \mathbf{a}\mathbf{s}' + \mathbf{e}'$
	$\mathbf{v} \leftarrow \mathbf{b}\mathbf{s}' + \mathbf{e}''$
	$\mathbf{r} \xleftarrow{\$} \text{HelpRec}(\mathbf{v})$
$\mathbf{v}' \leftarrow \mathbf{u}\mathbf{s}$	$\mathbf{v} \leftarrow \text{Rec}(\mathbf{v}, \mathbf{r})$
$\mathbf{v} \leftarrow \text{Rec}(\mathbf{v}', \mathbf{r})$	$\mu \leftarrow \text{SHA3-256}(\mathbf{v})$
$\mu \leftarrow \text{SHA3-256}(\mathbf{v})$	

Protocol 2: Our Scheme. For the definitions of HelpRec and Rec see Section 5. For the definition of encodings and the definition of Parse see Section 7.

most TLS implementations use fixed primes for Diffie-Hellman. One of the recommended mitigations by the authors of [1] is to avoid fixed primes.

Against all authority. Fortunately, all those pitfalls can be avoided by having the communicating parties generate a fresh \mathbf{a} at each instance of the protocol (as we propose). If in practice it turns out to be too expensive to generate \mathbf{a} for every connection, it is also possible to cache \mathbf{a} on the server side³ for, say a few hours without significantly weakening the protection against all-for-the-price-of-one attacks. Additionally, the performance impact of generating \mathbf{a} is reduced by sampling \mathbf{a} uniformly directly in NTT format (recalling that the NTT is a one-to-one map), and by transferring only a short 256-bit seed for \mathbf{a} (see Section 7).

A subtle question is to choose an appropriate primitive to generate a “random-looking” polynomial \mathbf{a} out of a short seed. For a security reduction, it seems to the authors that there is no way around the (non-programmable) random oracle model (ROM). It is argued in [38] that such a requirement is in practice an overkill, and that any pseudorandom generator (PRG) should also work. And while it is an interesting question how such a reasonable pseudo-random generator would interact with our lattice assumption, the cryptographic notion of a PRG *is not* helpful to argue security. Indeed, it is an easy exercise⁴ to build (under the NTRU assumption) a “backdoored” PRG that is, formally, a legitimate PRG, but that makes our scheme insecure.

³But recall that the secrets $\mathbf{s}, \mathbf{e}, \mathbf{s}', \mathbf{e}''$ have to be sampled fresh for every connection.

⁴Consider a secure PRG p , and parse its output $p(\text{seed})$ as two small polynomials (\mathbf{f}, \mathbf{g}) : an NTRU secret-key. Define $p'(\text{seed}) = \mathbf{g}\mathbf{f}^{-1} \bmod q$: under the decisional NTRU assumption, p' is still a secure PRG. Yet revealing the seed does reveal (\mathbf{f}, \mathbf{g}) and provides a backdoor as detailed above.

Instead, we prefer to base ourselves on a standard cryptographic hash-function, which is the typical choice of an “instantiation” of the ROM. As a suitable option we see Keccak [21], which has recently been standardized as SHA3 in FIPS-202 [75], and which offers extendable-output functions (XOF) named SHAKE. This avoids costly external iteration of a regular hash function and directly fits our needs.

We use SHAKE-128 for the generation of \mathbf{a} , which offers 128-bits of (post-quantum) security against collisions and preimage attacks. With only a small performance penalty we could have also chosen SHAKE-256, but we do not see any reason for such a choice, in particular because neither collisions nor preimages lead to an attack against the proposed scheme.

4 Choice of the error distribution

On non-Gaussian errors. In works like [22, 32, 85], a significant algorithmic effort is devoted to sample from a discrete Gaussian distribution to a rather high precision. In the following we argue that such effort is not necessary and motivate our choice of a centered binomial ψ_k as error distribution.

Indeed, we recall that the original worst-case to average-case reductions for LWE [83] and Ring-LWE [70] state hardness for *continuous Gaussian* distributions (and therefore also trivially apply to *rounded Gaussian*, which differ from discrete Gaussians). This also extends to discrete Gaussians [23] but such proofs are not necessarily intended for direct implementations. We recall that the use of discrete Gaussians (or other distributions with very high-precision sampling) is only crucial for signatures [68] and lattice trapdoors [43], to provide zero-knowledgeness.

The following Theorem states that choosing ψ_k as error distribution in Protocol 2 does not significantly decrease security compared to a rounded Gaussian distribution with the same standard deviation $\sigma = \sqrt{16/2}$.

Theorem 4.1 *Let ξ be the rounded Gaussian distribution of parameter $\sigma = \sqrt{8}$, that is, the distribution of $\lfloor \sqrt{8} \cdot x \rfloor$ where x follows the standard normal distribution. Let \mathcal{P} be the idealized version of Protocol 2, where the distribution ψ_{16} is replaced by ξ . If an (unbounded) algorithm, given as input the transcript of an instance of Protocol 2 succeeds in recovering the pre-hash key v with probability p , then it would also succeed against \mathcal{P} with probability at least*

$$q \geq p^{9/8}/26.$$

Proof See Appendix B.

As explained in Section 6, our choice of parameters leaves a comfortable margin to the targeted 128 bits of post-quantum security, which accommodates for the slight loss in security indicated by Theorem 4.1. Even more important from a practical point of view is that no known attack makes use of the difference in error distribution; what matters for attacks are entropy and standard deviation.

Simple implementation. We remark that sampling from the centered binomial distribution ψ_{16} is rather trivial in hardware and software, given the availability of a uniform binary source. Additionally, the implementation of this sampling algorithm is much easier to protect against timing attacks as no large tables or data-dependent branches are required (cf. to the issues caused by the table-based approach used in [22]).

5 Improved error-recovery mechanism

In most of the literature, Ring-LWE encryption allows to encrypt one bit per coordinate of the ciphertext. It is also well known how to encrypt multiple bits per coordinate by using a larger modulus-to-error ratio (and therefore decreasing the security for a fixed dimension n). However, in the context of exchanging a symmetric key (of, say, 256 bits), we end up having a message space larger than necessary and thus want to encrypt *one bit in multiple coordinates*.

In [82] Pöppelmann and Güneysu introduced a technique to encode one bit into two coordinates, and verified experimentally that it led to a better error tolerance. This allows to either increase the error and therefore improve the security of the resulting scheme or to decrease the probability of decryption failures. In this section we propose a generalization of this technique in dimension 4.

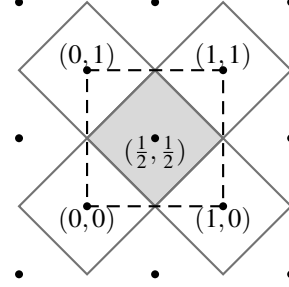


Figure 1: The lattice \tilde{D}_2 with Voronoi cells

We start with an intuitive description of the approach in 2 dimensions and then explain what changes in 4 dimensions. Appendices C and D give a thorough mathematical description together with a rigorous analysis.

Let us first assume that both client and server have the same vector $\mathbf{x} \in [0, 1)^2 \subset \mathbb{R}^2$ and want to map this vector to a single bit. Mapping polynomial coefficients from $\{0, \dots, q-1\}$ to $[0, 1)$ is easily accomplished through a division by q .

Now consider the lattice \tilde{D}_2 with basis $\{(0, 1), (\frac{1}{2}, \frac{1}{2})\}$. This lattice is a scaled version of the root lattice D_2 , specifically, $\tilde{D}_2 = \frac{1}{2} \cdot D_2$. Part of \tilde{D}_2 is depicted in Figure 1; lattice points are shown together with their Voronoi cells and the possible range of the vector \mathbf{x} is marked with dashed lines. Mapping \mathbf{x} to one bit is done by finding the closest-vector $v \in \tilde{D}_2$. If $v = (\frac{1}{2}, \frac{1}{2})$ (i.e., \mathbf{x} is in the grey Voronoi cell), then the output bit is 1; if $v \in \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ (i.e., \mathbf{x} is in a white Voronoi cell) then the output bit is 0.

This map may seem like a fairly complex way to map from a vector to a bit. However, recall that client and server only have a noisy version of \mathbf{x} , i.e., the client has a vector \mathbf{x}_c and the server has a vector \mathbf{x}_s . Those two vectors are close, but they are not the same and can be on different sides of a Voronoi cell border.

Error reconciliation. The approach described above now allows for an efficient solution to solve this agreement-from-noisy-data problem. The idea is that one of the two participants (in our case the client) sends as a *reconciliation vector* the difference of his vector \mathbf{x}_c and the center of its Voronoi cell (i.e., the point in the lattice). The server adds this difference vector to \mathbf{x}_s and thus moves away from the border towards the center of the correct Voronoi cell. Note that an eavesdropper does not learn anything from the reconciliation information: the client tells the difference to a lattice point, but not whether this is a lattice point producing a zero bit or a one bit.

This approach would require sending a full additional vector; we can reduce the amount of reconciliation in-

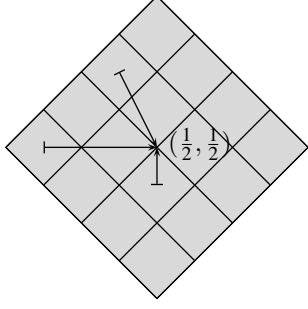


Figure 2: Splitting of the Voronoi cell of $(\frac{1}{2}, \frac{1}{2})$ into $2^{rd} = 16$ sub-cells, some with their corresponding difference vector to the center

formation through r -bit discretization. The idea is to split each Voronoi cell into 2^{dr} sub-cells and only send in which of those sub-cells the vector x_c is. Both participants then add the difference of the center of the sub-cell and the lattice point. This is illustrated for $r = 2$ and $d = 2$ in Figure 2.

Blurring the edges. Figure 1 may suggest that the probability of \mathbf{x} being in a white Voronoi cell is the same as for \mathbf{x} being in the grey Voronoi cell. This would be the case if \mathbf{x} actually followed a continuous uniform distribution. However, the coefficients of \mathbf{x} are discrete values in $\{0, \frac{1}{q}, \dots, \frac{q-1}{q}\}$ and with the protocol described so far, the bits of \mathbf{v} would have a small bias. The solution is to add, with probability $\frac{1}{2}$, the vector $(\frac{1}{2q}, \frac{1}{2q})$ to \mathbf{x} before running the error reconciliation. This has close to no effect for most values of \mathbf{x} , but, with probability $\frac{1}{2}$ moves \mathbf{x} to another Voronoi cell if it is very close to one side of a border. Appendix E gives a graphical intuition for this trick in two dimensions and with $q = 9$. The proof that it indeed removes all biases in the key is given in Lemma C.2.

From 2 to 4 dimensions. When moving from the 2-dimensional case considered above to the 4-dimensional case used in our protocol, not very much needs to change. The lattice \tilde{D}_2 becomes the lattice \tilde{D}_4 with basis $\mathbf{B} = (\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2, \mathbf{g})$, where \mathbf{u}_i are the canonical basis vectors of \mathbb{Z}^4 and $\mathbf{g}^t = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$. The lattice \tilde{D}_4 is a rotated and scaled version of the root lattice D_4 . The Voronoi cells of this lattice are no longer 2-dimensional “diamonds”, but 4-dimensional objects called icositetrachoron or 24-cells [64]. Determining in which cell a target point lies in is done using the closest vector algorithm $\text{CVP}_{\tilde{D}_4}$, and a simplified version of it, which we call Decode, gives the result modulo \mathbb{Z}^4 .

As in the 2-dimensional illustration in Figure 2, we are using 2-bit discretization; we are thus sending $r \cdot d = 8$ bits of reconciliation information per key bit.

Putting all of this together, we obtain the HelpRec function to compute the r -bit reconciliation information as

$$\text{HelpRec}(\mathbf{x}; b) = \text{CVP}_{\tilde{D}_4} \left(\frac{2^r}{q} (\mathbf{x} + b\mathbf{g}) \right) \bmod 2^r,$$

where $b \in \{0, 1\}$ is a uniformly chosen random bit. The corresponding function $\text{Rec}(\mathbf{x}, \mathbf{r}) = \text{Decode}(\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{B}\mathbf{r})$ computes one key bit from a vector \mathbf{x} with 4 coefficients in \mathbb{Z}_q and a reconciliation vector $\mathbf{r} \in \{0, 1, 2, 3\}^4$. The algorithms $\text{CVP}_{\tilde{D}_4}$ and Decode are listed as Algorithm 1 and Algorithm 2, respectively.

Algorithm 1 $\text{CVP}_{\tilde{D}_4}(\mathbf{x} \in \mathbb{R}^4)$

Ensure: An integer vector \mathbf{z} such that \mathbf{Bz} is a closest vector to \mathbf{x} : $\mathbf{x} - \mathbf{Bz} \in \mathcal{V}$

- 1: $\mathbf{v}_0 \leftarrow \lfloor \mathbf{x} \rfloor$
 - 2: $\mathbf{v}_1 \leftarrow \lfloor \mathbf{x} - \mathbf{g} \rfloor$
 - 3: $k \leftarrow (\|\mathbf{x} - \mathbf{v}_0\|_1 < 1) ? 0 : 1$
 - 4: $(v_0, v_1, v_2, v_3)^t \leftarrow \mathbf{v}_k$
 - 5: **return** $(v_0, v_1, v_2, k)^t + v_3 \cdot (-1, -1, -1, 2)^t$
-

Algorithm 2 $\text{Decode}(\mathbf{x} \in \mathbb{R}^4 / \mathbb{Z}^4)$

Ensure: A bit k such that $k\mathbf{g}$ is a closest vector to $\mathbf{x} + \mathbb{Z}^4$: $\mathbf{x} - k\mathbf{g} \in \mathcal{V} + \mathbb{Z}^4$

- 1: $\mathbf{v} = \mathbf{x} - \lfloor \mathbf{x} \rfloor$
 - 2: **return** 0 if $\|\mathbf{v}\|_1 \leq 1$ and 1 otherwise
-

Finally it remains to remark that even with this reconciliation mechanism client and server do not always agree on the same key. Appendix D provides a detailed analysis of the failure probability of the key agreement and shows that it is smaller than 2^{-60} .

6 Post-quantum security analysis

In [22] the authors chose Ring-LWE for a ring of rank $n = 1024$, while most previous instantiations of the Ring-LWE encryption scheme, like the ones in [32, 46, 66, 82], chose substantially smaller rank $n = 256$ or $n = 512$. It is argued that it is unclear if dimension 512 can offer post-quantum security. Yet, the concrete post-quantum security of LWE-based schemes has not been thoroughly studied, as far as we know. In this section we propose such a (very pessimistic) concrete analysis. In particular, our analysis reminds us that the security depends as much on q and its ratio with the error standard deviation ζ as it does on the dimension n . That means that our effort of optimizing the error recovery and its analysis not only improves efficiency but also offers superior security.

Security level over-shoot? With all our improvements, it would be possible to build a scheme with $n = 512$ (and $k = 24$, $q = 12289$) and to obtain security somewhat similar to the one of [22, 46], and therefore further improve efficiency. We call this variant JARJAR and details are provided in Appendix A. Nevertheless, as history showed us with RSA-512 [31], the standardization and deployment of a scheme awakens further cryptanalytic effort. In particular, NEWHOPE could withstand a dimension-halving attack in the line of [40, Sec 8.8.1] based on the Gentry-Szydlo algorithm [44, 63] or the sub-field approach of [2]. Note that so far, such attacks are only known for principal ideal lattices or NTRU lattices, and there are serious obstructions to extend them to Ring-LWE, but such precaution seems reasonable until lattice cryptanalysis stabilizes.

We provide the security and performance analysis of JARJAR in Appendix A mostly for comparison with other lower-security proposals. We strongly recommend NEWHOPE for any immediate applications, and advise against using JARJAR until concrete cryptanalysis of lattice-based cryptography is better understood.

6.1 Methodology: the core SVP hardness

We analyze the hardness of Ring-LWE as an LWE problem, since, so far, the best known attacks do not make use of the ring structure. There are many algorithms to consider in general (see the survey [3]), yet many of those are irrelevant for our parameter set. In particular, because there are only $m = n$ samples available one may rule out BKW types of attacks [56] and linearization attacks [6]. This essentially leaves us with two BKZ [28, 86] attacks, usually referred to as primal and dual attacks that we will briefly recall below.

The algorithm BKZ proceeds by reducing a lattice basis using an SVP oracle in a smaller dimension b . It is known [51] that the number of calls to that oracle remains polynomial, yet concretely evaluating the number of calls is rather painful, and this is subject to new heuristic ideas [27, 28]. We choose to ignore this polynomial factor, and rather evaluate only the *core SVP hardness*, that is the cost of *one call* to an SVP oracle in dimension b , which is clearly a pessimistic estimation (from the defender’s point of view).

6.2 Enumeration versus quantum sieve

Typical implementations [25, 28, 39] use an enumeration algorithm as this SVP oracle, yet this algorithm runs in super-exponential time. On the other hand, the sieve algorithms are known to run in exponential time, but are so far slower in practice for accessible dimensions $b \approx 130$. We choose the latter to predict the core hardness and will

argue that for the targeted dimension, enumerations are expected to be greatly slower than sieving.

Quantum sieve. A lot of recent work has pushed the efficiency of the original lattice sieve algorithms [72, 78], improving the heuristic complexity from $(4/3)^{b+o(b)} \approx 2^{0.415b}$ down to $\sqrt{3/2}^{b+o(b)} \approx 2^{0.292b}$ (see [12, 58]). The hidden sub-exponential factor is known to be much greater than one in practice, so again, estimating the cost ignoring this factor leaves us with a significant pessimistic margin.

Most of those algorithms have been shown [57, 59] to benefit from Grover’s quantum search algorithm, bringing the complexity down to $2^{0.265b}$. It is unclear if further improvements are to be expected, yet, because all those algorithms require classically building lists of size $\sqrt{4/3}^{b+o(b)} \approx 2^{0.2075b}$, it is very plausible that the best quantum SVP algorithm would run in time greater than $2^{0.2075b}$.

Irrelevance of enumeration for our analysis. In [28], predictions of the cost of solving SVP classically using the most sophisticated heuristic enumeration algorithms are given. For example, solving SVP in dimension 100 requires visiting about 2^{39} nodes, and 2^{134} nodes in dimension 250. Because this enumeration is a backtracking algorithm, it does benefit from the recent quasi-quadratic speedup [73], decreasing the quantum cost to about at least 2^{20} to 2^{67} operations as the dimension increases from 100 to 250.

On the other hand, our best-known attack bound $2^{0.265b}$ gives a cost of 2^{66} in dimension 250, and the best plausible attack bound $2^{0.2075b} \approx 2^{39}$. Because enumeration is super-exponential (both in theory and practice), its cost will be worse than our bounds in dimension larger than 250 and we may safely ignore this kind of algorithm.⁵

6.3 Primal attack

The primal attack consists of constructing a unique-SVP instance from the LWE problem and solving it using BKZ. We examine how large the block dimension b is required to be for BKZ to find the unique solution. Given the matrix LWE instance $(\mathbf{A}, \mathbf{b} = \mathbf{A}\mathbf{s} + \mathbf{e})$ one builds the lattice $\Lambda = \{\mathbf{x} \in \mathbb{Z}^{m+n+1} : (\mathbf{A} | -\mathbf{I}_m | -\mathbf{b})\mathbf{x} = \mathbf{0} \bmod q\}$ of dimension $d = m + n + 1$, volume q^m , and with a unique-SVP solution $\mathbf{v} = (\mathbf{s}, \mathbf{e}, 1)$ of norm $\lambda \approx \zeta \sqrt{n+m}$. Note that the number of used samples m may be chosen between 0 and $2n$ in our case and we numerically optimize this choice.

⁵The numbers are taken from the latest full version of [28] available at http://www.di.ens.fr/~ychen/research/Full_BKZ.pdf.

Success condition. We model the behavior of BKZ using the geometric series assumption (which is known to be optimistic from the attacker’s point of view), that finds a basis whose Gram-Schmidt norms are given by $\|\mathbf{b}_i^*\| = \delta^{d-2i-1} \cdot \text{Vol}(\Lambda)^{1/d}$ where $\delta = ((\pi b)^{1/b} \cdot b/2\pi e)^{1/2(b-1)}$ [3, 27]. The unique short vector \mathbf{v} will be detected if the projection of \mathbf{v} onto the vector space spanned by the last b Gram-Schmidt vectors is shorter than \mathbf{b}_{d-b}^* . Its projected norm is expected to be $\zeta\sqrt{b}$, that is the attack is successful if and only if

$$\zeta\sqrt{b} \leq \delta^{2b-d-1} \cdot q^{m/d}. \quad (1)$$

6.4 Dual attack

The dual attack consists of finding a short vector in the dual lattice $\mathbf{w} \in \Lambda' = \{(\mathbf{x}, \mathbf{y}) \in \mathbb{Z}^m \times \mathbb{Z}^n : \mathbf{A}'\mathbf{x} = \mathbf{y} \bmod q\}$. Assume we have found a vector (\mathbf{x}, \mathbf{y}) of length ℓ and compute $z = \mathbf{v}^t \cdot \mathbf{b} = \mathbf{v}^t \mathbf{A}\mathbf{s} + \mathbf{v}^t \mathbf{e} = \mathbf{w}^t \mathbf{s} + \mathbf{v}^t \mathbf{e} \bmod q$ which is distributed as a Gaussian of standard deviation $\ell\zeta$ if (\mathbf{A}, \mathbf{b}) is indeed an LWE sample (otherwise it is uniform mod q). Those two distributions have maximal variation distance bounded by⁶ $\varepsilon = 4\exp(-2\pi^2\tau^2)$ where $\tau = \ell\zeta/q$, that is, given such a vector of length ℓ one has an advantage ε against decision-LWE.

The length ℓ of a vector given by the BKZ algorithm is given by $\ell = \|\mathbf{b}_0\|$. Knowing that Λ' has dimension $d = m + n$ and volume q^n we get $\ell = \delta^{d-1}q^{n/d}$. Therefore, obtaining an ε -distinguisher requires running BKZ with block dimension b where

$$-2\pi^2\tau^2 \geq \ln(\varepsilon/4). \quad (2)$$

Note that small advantages ε are not relevant since the agreed key is hashed: an attacker needs an advantage of at least $1/2$ to significantly decrease the search space of the agreed key. He must therefore amplify his success probability by building about $1/\varepsilon^2$ many such short vectors. Because the sieve algorithms provide $2^{0.2075b}$ vectors, the attack must be repeated at least R times where

$$R = \max(1, 1/(2^{0.2075b}\varepsilon^2)).$$

This makes the conservative assumption that all the vectors provided by the Sieve algorithm are as short as the shortest one.

6.5 Security claims

According to our analysis, we claim that our proposed parameters offer at least (and quite likely with a large

⁶A preliminary version of this paper contained a bogus formula for ε leading to under-estimating the cost of the dual attack. Correcting this formula leads to better security claim, and almost similar cost for the primal and dual attacks.

Attack	m	b	Known Classical	Known Quantum	Best Plausible
BCNS proposal [22]: $q = 2^{32} - 1, n = 1024, \zeta = 3.192$					
Primal	1062	296	86	78	61
Dual	1055	296	86	78	61
NTRUENCRYPT [53]: $q = 2^{12}, n = 743, \zeta \approx \sqrt{2/3}$					
Primal	613	603	176	159	125
Dual	635	600	175	159	124
JARJAR: $q = 12289, n = 512, \zeta = \sqrt{12}$					
Primal	623	449	131	119	93
Dual	602	448	131	118	92
NEWHOPE: $q = 12289, n = 1024, \zeta = \sqrt{8}$					
Primal	1100	967	282	256	200
Dual	1099	962	281	255	199

Table 1: Core hardness of NEWHOPE and JARJAR and selected other proposals from the literature. The value b denotes the block dimension of BKZ, and m the number of used samples. Cost is given in \log_2 and is the smallest cost for all possible choices of m and b . Note that our estimation is very optimistic about the abilities of the attacker so that our result for the parameter set from [22] *does not* indicate that it can be broken with $\approx 2^{80}$ bit operations, given today’s state-of-the-art in cryptanalysis.

margin) a post-quantum security of 128 bits. The cost of the primal attack and dual attacks (estimated by our script `scripts/PQsecurity.py`) are given in Table 1. For comparison we also give a lower bound on the security of [22] and do notice a significantly improved security in our proposal. Yet, because of the numerous pessimistic assumption made in our analysis, we do not claim any quantum attacks reaching those bounds.

Most other RLWE proposals achieve considerably lower security than NEWHOPE; for example, the highest-security parameter set used for RLWE encryption in [46] is very similar to the parameters of JARJAR. The situation is different for NTRUENCRYPT, which has been instantiated with parameters that achieve about 128 bits of security according to our analysis⁷.

Specifically, we refer to NTRUENCRYPT with $n = 743$ as suggested in [53]. A possible advantage of NTRUENCRYPT compared to NEWHOPE is somewhat smaller message sizes, however, this advantage becomes very small when scaling parameters to achieve a similar security margin as NEWHOPE. The large downside of using NTRUENCRYPT for ephemeral key exchange is the cost for key generation. The implemen-

⁷For comparison we view the NTRU key-recovery as an homogeneous Ring-LWE instance. We do not take into account the combinatorial vulnerabilities [55] induced by the fact that secrets are ternary. We note that NTRU is a potentially a weaker problem than Ring-LWE: it is in principle subject to a subfield-lattice attack [2], but the parameters proposed for NTRUENCRYPT are immune.

tation of NTRUENCRYPT with $n = 743$ in eBACS [19] takes about an order of magnitude longer for key generation alone than NEWHOPE takes in total. Also, unlike our NEWHOPE software, this NTRUENCRYPT software is not protected against timing attacks; adding such protection would presumably incur a significant overhead.

7 Implementation

In this section we provide details on the encodings of messages and describe our portable reference implementation written in C, as well as an optimized implementation targeting architectures with AVX vector instructions.

7.1 Encodings and generation of \mathbf{a}

The key-exchange protocol described in Protocol 1 and also our protocol as described in Protocol 2 exchange messages that contain mathematical objects (in particular, polynomials in \mathcal{R}_q). Implementations of these protocols need to exchange messages in terms of byte arrays. As we will describe in the following, the choice of encodings of polynomials to byte arrays has a serious impact on performance. We use an encoding of messages that is particularly well-suited for implementations that make use of quasi-linear NTT-based polynomial multiplication.

Definition of NTT and NTT^{-1} . The NTT is a tool commonly used in implementations of ideal lattice-based cryptography [32, 46, 66, 82]. For some background on the NTT and the description of fast software implementations we refer to [50, 71]. In general, fast quasi-logarithmic algorithms exist for the computation of the NTT and a polynomial multiplication can be performed by computing $\mathbf{c} = \text{NTT}^{-1}(\text{NTT}(\mathbf{a}) \circ \text{NTT}(\mathbf{b}))$ for $\mathbf{a}, \mathbf{b}, \mathbf{c} \in \mathcal{R}$. An NTT targeting ideal lattices defined in $\mathcal{R}_q = \mathbb{Z}_q[X]/(X^n + 1)$ can be implemented very efficiently if n is a power of two and q is a prime for which it holds that $q \equiv 1 \pmod{2n}$. This way a primitive n -th root of unity ω and its square root γ exist. By multiplying coefficient-wise by powers of $\gamma = \sqrt{\omega} \pmod{q}$ before the NTT computation and after the reverse transformation by powers of γ^{-1} , no zero padding is required and an n -point NTT can be used to transform a polynomial with n coefficients.

For a polynomial $\mathbf{g} = \sum_{i=0}^{1023} g_i X^i \in \mathcal{R}_q$ we define

$$\text{NTT}(\mathbf{g}) = \hat{\mathbf{g}} = \sum_{i=0}^{1023} \hat{g}_i X^i, \text{ with}$$

$$\hat{g}_i = \sum_{j=0}^{1023} \gamma^j g_j \omega^{ij},$$

where we fix the n -th primitive root of unity to $\omega = 49$ and thus $\gamma = \sqrt{\omega} = 7$. Note that in our implementation we use an in-place NTT algorithm which requires bit-reversal operations. As an optimization, our implementations skips these bit-reversals for the forward transformation as all inputs are only random noise. This optimization is transparent to the protocol and for simplicity omitted in the description here.

The function NTT^{-1} is the inverse of the function NTT. The computation of NTT^{-1} is essentially the same as the computation of NTT, except that it uses $\omega^{-1} \pmod{q} = 1254$, multiplies by powers of $\gamma^{-1} \pmod{q} = 8778$ after the summation, and also multiplies each coefficient by the scalar $n^{-1} \pmod{q} = 12277$ so that

$$\text{NTT}^{-1}(\hat{\mathbf{g}}) = \mathbf{g} = \sum_{i=0}^{1023} g_i X^i, \text{ with}$$

$$g_i = n^{-1} \gamma^{-i} \sum_{j=0}^{1023} \hat{g}_j \omega^{-ij}.$$

The inputs to NTT^{-1} are *not* just random noise, so inside NTT^{-1} our software has to perform the initial bit reversal, making NTT^{-1} slightly more costly than NTT.

Definition of Parse. The public parameter \mathbf{a} is generated from a 256-bit seed through the extendable-output function SHAKE-128 [75, Sec. 6.2]. The output of SHAKE-128 is considered as an array of 16-bit, unsigned, little-endian integers. Each of those integers is used as a coefficient of \mathbf{a} if it is smaller than $5q$ and rejected otherwise. The first such 16-bit integer is used as the coefficient of X^0 , the next one as coefficient of X^1 and so on. Earlier versions of this paper described a slightly different way of rejection sampling for coefficients of \mathbf{a} . The more efficient approach adopted in this final version was suggested independently by Gueron and Schlieker in [49] and by Yawning Angel in [5]. However, note that a reduction modulo q of the coefficients of \mathbf{a} as described in [49] and [5] is not necessary; both our implementations can handle coefficients of \mathbf{a} in $\{0, \dots, 5q - 1\}$.

Due to a small probability of rejections, the amount of output required from SHAKE-128 depends on the seed – what is required is $n = 1024$ coefficients that are smaller than $5q$. The minimal amount of output is thus 2 KB; the average amount is ≈ 2184.5 bytes. The resulting polynomial \mathbf{a} (denoted as $\hat{\mathbf{a}}$) is considered to be in NTT domain. This is possible because the NTT transforms uniform noise to uniform noise.

Using a variable amount of output from SHAKE-128 leaks information about \mathbf{a} through timing information. This is not a problem for most applications, since \mathbf{a} is public. As pointed out by Burdges in [24], such a timing leak of public information can be a problem when deploying NEWHOPE in anonymity networks like Tor.

Appendix F describes an alternative approach for Parse, which is slightly more complex and slightly slower, but does not leak any timing information about \mathbf{a} .

The message format of $(\mathbf{b}, \text{seed})$ and (\mathbf{u}, \mathbf{r}) . With the definition of the NTT, we can now define the format of the exchanged messages. In both $(\mathbf{b}, \text{seed})$ and (\mathbf{u}, \mathbf{r}) the polynomial is transmitted in the NTT domain (as in works like [82, 85]). Polynomials are encoded as an array of 1792 bytes, in a compressed little-endian format. The encoding of seed is straight-forward as an array of 32 bytes, which is simply concatenated with the encoding of \mathbf{b} . Also the encoding of \mathbf{r} is fairly straight-forward: it packs four 2-bit coefficients into one byte for a total of 256 bytes, which are again simply concatenated with the encoding of \mathbf{u} . We denote these encodings to byte arrays as `encodeA` and `encodeB` and their inverses as `decodeA` and `decodeB`. For a description of our key-exchange protocol including encodings and with explicit NTT and NTT⁻¹ transformations, see Protocol 3.

7.2 Portable C implementation

This paper is accompanied by a C reference implementation described in this section and an optimized implementation for Intel and AMD CPUs described in the next section. The main emphasis in the C reference implementation is on simplicity and portability. It does not use any floating-point arithmetic and outside the Keccak (SHA3-256 and SHAKE-128) implementation only needs 16-bit and 32-bit integer arithmetic. In particular, the error-recovery mechanism described in Section 5 is implemented with fixed-point (i.e., integer-) arithmetic. Furthermore, the C reference implementation does not make use of the division operator (`/`) and the modulo operator (`%`). The focus on simplicity and portability does not mean that the implementation is not optimized at all. On the contrary, we use it to illustrate various optimization techniques that are helpful to speed up the key exchange and are also of independent interest for implementers of other ideal-lattice-based schemes.

NTT optimizations. All polynomial coefficients are represented as unsigned 16-bit integers. Our in-place NTT implementation transforms from bit-reversed to natural order using Gentleman-Sande butterfly operations [29, 41]. One would usually expect that each NTT is preceded by a bit-reversal, but all inputs to NTT are noise polynomials that we can simply consider as being already bit-reversed; as explained earlier, the NTT⁻¹ operation still involves a bit-reversal. The core of the NTT and NTT⁻¹ operation consists of 10 layers of transformations, each consisting of 512 butterfly operations of the form described in Listing 2.

Montgomery arithmetic and lazy reductions. The performance of operations on polynomials is largely determined by the performance of NTT and NTT⁻¹. The main computational bottleneck of those operations are 5120 butterfly operations, each consisting of one addition, one subtraction and one multiplication by a precomputed constant. Those operations are in \mathbb{Z}_q ; recall that q is a 14-bit prime. To speed up the modular-arithmetic operations, we store all precomputed constants in Montgomery representation [74] with $R = 2^{18}$, i.e., instead of storing ω^i , we store $2^{18}\omega^i \pmod{q}$. After a multiplication of a coefficient g by some constant $2^{18}\omega^i$, we can then reduce the result r to $g\omega^i \pmod{q}$ with the fast Montgomery reduction approach. In fact, we do not always fully reduce modulo q , it is sufficient if the result of the reduction has at most 14 bits. The fast Montgomery reduction routine given in Listing 1a computes such a reduction to a 14-bit integer for any unsigned 32-bit integer in $\{0, \dots, 2^{32} - q(R - 1) - 1\}$. Note that the specific implementation does not work for *any* 32-bit integer; for example, for the input $2^{32} - q(R - 1) = 1073491969$ the addition `a=a+u` causes an overflow and the function returns 0 instead of the correct result 4095. In the following we establish that this is not a problem for our software.

Aside from reductions after multiplication, we also need modular reductions after addition. For this task we use the “short Barrett reduction” [10] detailed in Listing 1b. Again, this routine does not fully reduce modulo q , but reduces any 16-bit unsigned integer to an integer of at most 14 bits which is congruent modulo q .

In the context of the NTT and NTT⁻¹, we make sure that inputs have coefficients of at most 14 bits. This allows us to avoid Barrett reductions after addition on every second level, because coefficients grow by at most one bit per level and the short Barrett reduction can handle 16-bit inputs. Let us turn our focus to the input of the Montgomery reduction (see Listing 2). Before subtracting `a[j+d]` from `t` we need to add a multiple of q to avoid unsigned underflow. Coefficients never grow larger than 15 bits and $3 \cdot q = 36867 > 2^{15}$, so adding $3 \cdot q$ is sufficient. An upper bound on the expression `((uint32_t)t + 3*12289 - a[j+d])` is obtained if `t` is $2^{15} - 1$ and `a[j+d]` is zero; we thus obtain $2^{15} + 3 \cdot q = 69634$. All precomputed constants are in $\{0, \dots, q - 1\}$, so the expression `(w * ((uint32_t)t + 3*12289 - a[j+d]))`, the input to the Montgomery reduction, is at most $69634 \cdot (q - 1) = 855662592$ and thus safely below the maximum input that the Montgomery reduction can handle.

Fast random sampling. As a first step before performing any operations on polynomials, both Alice and Bob need to expand the seed to the polynomial \mathbf{a} using SHAKE-128. The implementation we use is based on the “simple” implementation by Van Keer for the Kecc-

Parameters: $q = 12289 < 2^{14}$, $n = 1024$	
Error distribution: ψ_{16}^n	
Alice (server)	Bob (client)
$seed \xleftarrow{\$} \{0, \dots, 255\}^{32}$	
$\hat{\mathbf{a}} \leftarrow \text{Parse}(\text{SHAKE-128}(seed))$	
$\mathbf{s}, \mathbf{e} \xleftarrow{\$} \psi_{16}^n$	$\mathbf{s}', \mathbf{e}', \mathbf{e}'' \xleftarrow{\$} \psi_{16}^n$
$\hat{\mathbf{s}} \leftarrow \text{NTT}(\mathbf{s})$	
$\hat{\mathbf{b}} \leftarrow \hat{\mathbf{a}} \circ \hat{\mathbf{s}} + \text{NTT}(\mathbf{e})$	$(\hat{\mathbf{b}}, seed) \leftarrow \text{decodeA}(m_a)$
	$\xrightarrow[1824 \text{ Bytes}]{m_a = \text{encodeA}(seed, \hat{\mathbf{b}})}$
	$\hat{\mathbf{a}} \leftarrow \text{Parse}(\text{SHAKE-128}(seed))$
	$\hat{\mathbf{t}} \leftarrow \text{NTT}(\mathbf{s}')$
	$\hat{\mathbf{u}} \leftarrow \hat{\mathbf{a}} \circ \hat{\mathbf{t}} + \text{NTT}(\mathbf{e}')$
	$\mathbf{v} \leftarrow \text{NTT}^{-1}(\hat{\mathbf{b}} \circ \hat{\mathbf{t}}) + \mathbf{e}''$
	$\xleftarrow[2048 \text{ Bytes}]{m_b = \text{encodeB}(\hat{\mathbf{u}}, \mathbf{r})}$
$(\hat{\mathbf{u}}, \mathbf{r}) \leftarrow \text{decodeB}(m_b)$	$\mathbf{r} \xleftarrow{\$} \text{HelpRec}(\mathbf{v})$
$\mathbf{v}' \leftarrow \text{NTT}^{-1}(\hat{\mathbf{u}} \circ \hat{\mathbf{s}})$	$\mathbf{v} \leftarrow \text{Rec}(\mathbf{v}, \mathbf{r})$
$\mathbf{v} \leftarrow \text{Rec}(\mathbf{v}', \mathbf{r})$	$\mu \leftarrow \text{SHA3-256}(\mathbf{v})$
$\mu \leftarrow \text{SHA3-256}(\mathbf{v})$	

Protocol 3: Our proposed protocol including NTT and NTT^{-1} computations and sizes of exchanged messages; \circ denotes pointwise multiplication; elements in NTT domain are denoted with a hat ($\hat{\cdot}$)

Listing 1 Reduction routines used in the reference implementation.

(a) Montgomery reduction ($R = 2^{18}$).

```
uint16_t mred(uint32_t a) {
    uint32_t u;
    u = (a * 12287);
    u &= ((1 << 18) - 1);
    a += u * 12289;
    return a >> 18;
}
```

(b) Short Barrett reduction.

```
uint16_t bred(uint16_t a) {
    uint32_t u;
    u = ((uint32_t) a * 5) >> 16;
    a -= u * 12289;
    return a;
}
```

caik permutation and slightly modified code taken from the “TweetFIPS202” implementation [20] for everything else.

The sampling of centered binomial noise polynomials is based on a fast PRG with a random seed from `/dev/urandom` followed by a quick summation of 16-bit chunks of the PRG output. Note that the choice of the PRG is a purely local choice that every user can pick independently based on the target hardware architecture and based on routines that are available anyway (for example, for symmetric encryption following the key exchange). Our C reference implementation uses ChaCha20 [14], which is fast, trivially protected against timing attacks, and is already in use by many TLS clients

Listing 2 The Gentleman-Sande butterfly inside odd levels of our NTT computation. All $a[j]$ and W are of type `uint16_t`.

```
W = omega[jTwiddle++];
t = a[j];
a[j] = bred(t + a[j+d]);
a[j+d] = mred(W * ((uint32_t)t + 3*12289 - a[j+d]));
```

and servers [60, 61].

7.3 Optimized AVX2 implementation

Intel processors since the “Sandy Bridge” generation support Advanced Vector Extensions (AVX) that operate on vectors of 8 single-precision or 4 double-precision floating-point values in parallel. With the introduction of the “Haswell” generation of CPUs, this support was extended also to 256-bit vectors of integers of various sizes (AVX2). It is not surprising that the enormous computational power of these vector instructions has been used before to implement very high-speed crypto (see, for example, [16, 18, 47]) and also our optimized reference implementation targeting Intel Haswell processors uses those instructions to speed up multiple components of the key exchange.

NTT optimizations. The AVX instruction set has been used before to speed up the computation of lattice-based cryptography, and in particular the number-theoretic transform. Most notably, Güneysu, Oder, Pöppelmann

and Schwabe achieve a performance of only 4480 cycles for a dimension-512 NTT on Intel Sandy Bridge [50]. For arithmetic modulo a 23-bit prime, they represent coefficients as double-precision integers.

We experimented with multiple different approaches to speed up the NTT in AVX. For example, we vectorized the Montgomery arithmetic approach of our C reference implementation and also adapted it to a 32-bit-signed-integer approach. In the end it turned out that floating-point arithmetic beats all of those more sophisticated approaches, so we are now using an approach that is very similar to the approach in [50]. One computation of a dimension-1024 NTT takes 8448 cycles, unlike the numbers in [50] this does include multiplication by the powers of γ and unlike the numbers in [50], this excludes a bit-reversal.

Fast sampling. Intel Haswell processors support the AES-NI instruction set and for the local choice of noise sampling it is obvious to use those. More specifically, we use the public-domain implementation of AES-256 in counter mode written by Dolbeau, which is included in the SUPERCOP benchmarking framework [19]. Transformation from uniform noise to the centered binomial is optimized in AVX2 vector instructions operating on vectors of bytes and 16-bit integers.

For the computation of SHAKE-128 we use the same code as in the C reference implementation. One might expect that architecture-specific optimizations (for example, using AVX instructions) are able to offer significant speedups, but the benchmarks of the eBACS project [19] indicate that on Intel Haswell, the fastest implementation is the “simple” implementation by Van Keer that our C reference implementation is based on. The reasons that vector instructions are not very helpful for speeding up SHAKE (or, more generally, Keccak) are the inherently sequential nature and the 5×5 dimension of the state matrix that makes internal vectorization hard.

Error recovery. The 32-bit integer arithmetic used by the C reference implementation for HelpRec and Rec is trivially 8-way parallelized with AVX2 instructions. With this vectorization, the cost for HelpRec is only 3404 cycles, the cost for Rec is only 2804 cycles.

8 Benchmarks and comparison

In the following we present benchmark results of our software. All benchmark results reported in Table 2 were obtained on an Intel Core i7-4770K (Haswell) running at 3491.953 MHz with Turbo Boost and Hyperthreading disabled. We compiled our C reference implementation with gcc-4.9.2 and flags `-O3 -fomit-frame-pointer -march=corei7-avx -msse2avx`. We compiled our

optimized AVX implementation with clang-3.5 and flags `-O3 -fomit-frame-pointer -march=native`.

As described in Section 7, the sampling of \mathbf{a} is not running in constant time; we report the median running time and (in parentheses) the average running time for this generation, the server-side key-pair generation and client-side shared-key computation; both over 1000 runs. For all other routines we report the median of 1000 runs. We built the software from [22] on the same machine as ours and—like the authors of [22]—used `openssl speed` for benchmarking their software and converted the reported results to approximate cycle counts as given in Table 2.

Comparison with BCNS and RSA/ECDH. As previously mentioned, the BCNS implementation [22] also uses the dimension $n = 1024$ but the larger modulus $q = 2^{32} - 1$ and the Gaussian error distribution with Gaussian parameter $\sigma = 8/\sqrt{2\pi} = 3.192$. When the authors of BCNS integrated their implementation into SSL it only incurred a slowdown by a factor of 1.27 compared to ECDH when using ECDSA signatures and a factor of 1.08 when using RSA signatures with respect to the number of connections that could be handled by the server. As a reference, the reported cycle counts in [22] for a `nistp256` ECDH on the client side are 2 160 000 cycles (0.8 ms @2.77 GHz) and on the server side 3 221 288 cycles (1.4 ms @2.33 GHz). These numbers are obviously not state of the art for ECDH software. Even on the `nistp256` curve, which is known to be a far-from-optimal choice, it is possible to achieve cycle counts of less than 300 000 cycles for a variable-basepoint scalar multiplication on an Intel Haswell [48]. Also OpenSSL optionally includes fast software for `nistp256` ECDH by Käsper and Langley and we assume that the authors of [22] omitted enabling it. Compared to BCNS, our C implementation is more than 8 times faster and our AVX implementation even achieves a speedup factor of more than 27. At this performance it is in the same ballpark as state-of-the-art ECDH software, even when TLS switches to faster 128-bit secure ECDH key exchange based on Curve25519 [13], as recently specified in RFC 7748 [62].

In comparison to the BCNS proposal we see a large performance advantage from switching to the binomial error distribution. The BCNS software uses a large pre-computed table to sample from a discrete Gaussian distribution with a high precision. This approach takes 1042700 cycles to sample one polynomial in constant time. Our C implementation requires only 32684 cycles to sample from the binomial distribution. Another factor is that we use the NTT in combination with a smaller modulus. Polynomial multiplication in [22] is using Nussbaumer’s symbolic approach based on recursive negacyclic convolutions [79]. The implemen-

Table 2: Intel Haswell cycle counts of our proposal as compared to the BCNS proposal from [22].

	BCNS [22]	Ours (C ref)	Ours (AVX2)
Generation of a		43 440 ^a (43 607) ^a	37 470 ^a (36 863) ^a
NTT		55 360	8 448
NTT ⁻¹		59 864 ^b	9 464 ^b
Sampling of a noise polynomial		32 684 ^c	5 900 ^c
HelpRec		14 608	3 404
Rec		10 092	2 804
Key generation (server)	≈ 2 477 958	258 246 (258 965)	88 920 (89 079)
Key gen + shared key (client)	≈ 3 995 977	384 994 (385 146)	110 986 (111 169)
Shared key (server)	≈ 481 937	86 280	19 422

^a Includes reading a seed from `/dev/urandom`

^b Includes one bit reversal

^c Excludes reading a seed from `/dev/urandom`, which is shared across multiple calls to the noise generation

tation in [22] only achieves a performance of 342 800 cycles for a constant-time multiplication. Additionally, the authors of [22] did not perform pre-transformation of constants (e.g., **a**) or transmission of coefficients in FFT/Nussbaumer representation.

Follow-Up Work. We would like to refer the reader to follow-up work in which improvements to NEWHOPE and its implementation were proposed based on a preprint version of this work [4]. In [49] Gueron and Schlieker introduce faster pseudorandom bytes generation by changing the underlying functions, a method to decrease the rejection rate during sampling (see Section 7.1), and a vectorization of the sampling step. Longa and Naehrig [67] optimize the NTT and present new modular reduction techniques and are able to achieve a speedup of factor-1.90 for the C implementation and a factor-1.25 for the AVX implementation compared to the preprint [4] (note that this version has updated numbers). An alternative NTRU-based proposal and implementation of a lattice-based public-key encryption scheme, which could also be used for key exchange, is given by Bernstein, Chuengsatiansup, Lange, and van Vredendaal in [17], but we leave a detailed comparison to future work. An efficient authenticated lattice-based key exchange scheme has recently been proposed by del Pino, Lyubashevsky, and Pointcheval in [33].

References

[1] ADRIAN, D., BHARGAVAN, K., DURUMERIC, Z., GAUDRY, P., GREEN, M., HALDERMAN, J. A., HENINGER, N., SPRINGALL, D., THOMÉ, E., VALENTA, L., VANDERSLOOT, B., WUSTROW, E., BÉGUELIN, S. Z., AND ZIMMERMANN, P. Imperfect forward secrecy: How Diffie-Hellman fails in practice.

In *CCS '15 Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security* (2015), ACM, pp. 5–17. <https://weakdh.org/>. 4, 5

- [2] ALBRECHT, M., BAI, S., AND DUCAS, L. A subfield lattice attack on overstretched NTRU assumptions. *IACR Cryptology ePrint Archive report 2016/127*, 2016. <http://eprint.iacr.org/2016/127>. 8, 9
- [3] ALBRECHT, M. R., PLAYER, R., AND SCOTT, S. On the concrete hardness of learning with errors. *IACR Cryptology ePrint Archive report 2015/046*, 2015. <http://eprint.iacr.org/2015/046/>. 8, 9
- [4] ALKIM, E., DUCAS, L., PÖPELMANN, T., AND SCHWABE, P. Post-quantum key exchange – a new hope. *IACR Cryptology ePrint Archive report 2015/1092*, 2015. <https://eprint.iacr.org/2015/1092/20160329:201913>. 14
- [5] ANGEL, Y. Post-quantum secure hybrid handshake based on NewHope. Posting to the tor-dev mailing list, 2016. <https://lists.torproject.org/pipermail/tor-dev/2016-May/010896.html>. 10
- [6] ARORA, S., AND GE, R. New algorithms for learning in presence of errors. In *Automata, Languages and Programming* (2011), L. Aceto, M. Henzinger, and J. Sgall, Eds., vol. 6755 of *LNCS*, Springer, pp. 403–415. <https://www.cs.duke.edu/~rongge/LPSN.pdf>. 8, 18
- [7] BABAI, L. On Lovász’ lattice reduction and the nearest lattice point problem. *Combinatorica* 6, 1 (1986), 1–13. <http://www.csie.nuk.edu.tw/~cychen/Lattices/0n%20lovasz%20lattice%20reduction%20and%20the%20nearest%20lattice%20point%20problem.pdf>. 4
- [8] BAI, S., AND GALBRAITH, S. D. An improved compression technique for signatures based on learning with errors. In *Topics in Cryptology – CT-RSA 2014* (2014), J. Benaloh, Ed., vol. 8366 of *LNCS*, Springer, pp. 28–47. <https://eprint.iacr.org/2013/838/>. 1
- [9] BAI, S., LANGLOIS, A., LEPOINT, T., STEHLÉ, D., AND STEINFELD, R. Improved security proofs in lattice-based cryptography: using the Rényi divergence rather than the statistical distance. In *Advances in Cryptology – ASIACRYPT 2015*

- (2015), T. Iwata and J. H. Cheon, Eds., LNCS, Springer. <http://eprint.iacr.org/2015/483/>. 18
- [10] BARRETT, P. Implementing the Rivest Shamir and Adleman public key encryption algorithm on a standard digital signal processor. In *Advances in Cryptology – CRYPTO '86* (1987), A. M. Odlyzko, Ed., vol. 263 of *Lecture Notes in Computer Science*, Springer-Verlag Berlin Heidelberg, pp. 311–323. 11
- [11] BATCHER, K. E. Sorting networks and their applications. In *AFIPS conference proceedings, volume 32: 1968 Spring Joint Computer Conference* (1968), Thompson Book Company, pp. 307–314. <http://www.cs.kent.edu/~batcher/conf.html>. 22
- [12] BECKER, A., DUCAS, L., GAMA, N., AND LAARHOVEN, T. New directions in nearest neighbor searching with applications to lattice sieving. In *SODA '16 Proceedings of the twenty-seventh annual ACM-SIAM symposium on Discrete Algorithms* (2016 (to appear)), SIAM. 8
- [13] BERNSTEIN, D. J. Curve25519: new Diffie-Hellman speed records. In *Public Key Cryptography – PKC 2006* (2006), M. Yung, Y. Dodis, A. Kiayias, and T. Malkin, Eds., vol. 3958 of LNCS, Springer, pp. 207–228. <http://cr.yo.to/papers.html#curve25519>. 13
- [14] BERNSTEIN, D. J. ChaCha, a variant of Salsa20. In *Workshop Record of SASC 2008: The State of the Art of Stream Ciphers* (2008). <http://cr.yo.to/papers.html#chacha>. 12
- [15] BERNSTEIN, D. J., CHOU, T., CHUENGSATIANSUP, C., HÜLSING, A., LANGE, T., NIEDERHAGEN, R., AND VAN VREDENDAAL, C. How to manipulate curve standards: a white paper for the black hat. IACR Cryptology ePrint Archive report 2014/571, 2014. <http://eprint.iacr.org/2014/571/>. 4
- [16] BERNSTEIN, D. J., CHUENGSATIANSUP, C., LANGE, T., AND SCHWABE, P. Kummer strikes back: new DH speed records. In *Advances in Cryptology – EUROCRYPT 2015* (2014), T. Iwata and P. Sarkar, Eds., vol. 8873 of LNCS, Springer, pp. 317–337. full version: <http://cryptojedi.org/papers/#kummer>. 12
- [17] BERNSTEIN, D. J., CHUENGSATIANSUP, C., LANGE, T., AND VAN VREDENDAAL, C. NTRU Prime. IACR Cryptology ePrint Archive report 2016/461, 2016. <https://eprint.iacr.org/2016/461>. 14
- [18] BERNSTEIN, D. J., HOPWOOD, D., HÜLSING, A., LANGE, T., NIEDERHAGEN, R., PAPACHRISTODOULOU, L., SCHNEIDER, M., SCHWABE, P., AND WILCOX-O'HEARN, Z. SPHINCS: practical stateless hash-based signatures. In *Advances in Cryptology – EUROCRYPT 2015* (2015), E. Oswald and M. Fischlin, Eds., vol. 9056 of LNCS, Springer, pp. 368–397. <https://cryptojedi.org/papers/#sphincs>. 3, 12
- [19] BERNSTEIN, D. J., AND LANGE, T. eBACS: ECRYPT benchmarking of cryptographic systems. <http://bench.cr.yo.to> (accessed 2015-10-07). 10, 13
- [20] BERNSTEIN, D. J., SCHWABE, P., AND ASSCHE, G. V. Tweetable FIPS 202, 2015. <http://keccak.noekeon.org/tweetfips202.html> (accessed 2016-03-21). 12
- [21] BERTONI, G., DAEMEN, J., PEETERS, M., AND ASSCHE, G. V. Keccak. In *Advances in Cryptology – EUROCRYPT 2013* (2013), T. Johansson and P. Q. Nguyen, Eds., vol. 7881 of LNCS, Springer, pp. 313–314. 5
- [22] BOS, J. W., COSTELLO, C., NAEHRIG, M., AND STEBILA, D. Post-quantum key exchange for the TLS protocol from the ring learning with errors problem. In *2015 IEEE Symposium on Security and Privacy* (2015), pp. 553–570. <http://eprint.iacr.org/2014/599>. 1, 2, 3, 4, 5, 6, 7, 8, 9, 13, 14
- [23] BRAKERSKI, Z., LANGLOIS, A., PEIKERT, C., REGEV, O., AND STEHLÉ, D. Classical hardness of learning with errors. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing* (2013), ACM, pp. 575–584. <http://arxiv.org/pdf/1306.0281>. 5
- [24] BURDGES, J. Post-quantum secure hybrid handshake based on NewHope. Posting to the tor-dev mailing list, 2016. <https://lists.torproject.org/pipermail/tor-dev/2016-May/010886.html>. 10
- [25] CADÉ, D., PUJOL, X., AND STEHLÉ, D. fp11 4.0.4, 2013. <https://github.com/dstehle/fp11> (accessed 2015-10-13). 8
- [26] CHECKOWAY, S., FREDRIKSON, M., NIEDERHAGEN, R., EVERSPAUGH, A., GREEN, M., LANGE, T., RISTENPART, T., BERNSTEIN, D. J., MASKIEWICZ, J., AND SHACHAM, H. On the practical exploitability of Dual EC in TLS implementations. In *Proceedings of the 23rd USENIX security symposium* (2014). <https://projectbullrun.org/dual-ec/index.html>. 4
- [27] CHEN, Y. *Lattice reduction and concrete security of fully homomorphic encryption*. PhD thesis, l'Université Paris Diderot, 2013. Available at <http://www.di.ens.fr/~ychen/research/these.pdf>. 8, 9
- [28] CHEN, Y., AND NGUYEN, P. Q. BKZ 2.0: Better lattice security estimates. In *Advances in Cryptology – ASIACRYPT 2011*, D. H. Lee and X. Wang, Eds., vol. 7073 of LNCS, Springer, 2011, pp. 1–20. <http://www.iacr.org/archive/asiacrypt2011/70730001/70730001.pdf>. 8
- [29] CHU, E., AND GEORGE, A. *Inside the FFT Black Box Serial and Parallel Fast Fourier Transform Algorithms*. CRC Press, Boca Raton, FL, USA, 2000. 11
- [30] CONWAY, J., AND SLOANE, N. J. A. *Sphere packings, lattices and groups*, 3rd ed., vol. 290 of *Grundlehren der mathematischen Wissenschaften*. Springer Science & Business Media, 1999. 18
- [31] COWIE, J., DODSON, B., ELKENBRACHT-HUIZING, R. M., LENSTRA, A. K., MONTGOMERY, P. L., AND ZAYER, J. A world wide number field sieve factoring record: on to 512 bits. In *Advances in Cryptology – ASIACRYPT'96* (1996), K. Kim and T. Matsumoto, Eds., vol. 1163 of LNCS, Springer, pp. 382–394. <http://oai.cwi.nl/oai/asset/1940/1940A.pdf>. 8
- [32] DE CLERCQ, R., ROY, S. S., VERCAUTEREN, F., AND VERBAUWHEDE, I. Efficient software implementation of ring-LWE encryption. In *Design, Automation & Test in Europe Conference & Exhibition, DATE 2015* (2015), EDA Consortium, pp. 339–344. <http://eprint.iacr.org/2014/725>. 1, 3, 5, 7, 10
- [33] DEL PINO, R., LYUBASHEVSKY, V., AND POINTCHEVAL, D. The whole is less than the sum of its parts: Constructing more efficient lattice-based AKEs. IACR Cryptology ePrint Archive report 2016/435, 2016. <https://eprint.iacr.org/2016/435>. 14
- [34] DING, J., XIE, X., AND LIN, X. A simple provably secure key exchange scheme based on the learning with errors problem. IACR Cryptology ePrint Archive report 2012/688, 2012. <http://eprint.iacr.org/2012/688>. 1
- [35] DUCAS, L., DURMUS, A., LEPOINT, T., AND LYUBASHEVSKY, V. Lattice signatures and bimodal Gaussians. In *Advances in Cryptology – CRYPTO 2013* (2013), R. Canetti and J. A. Garay, Eds., vol. 8042 of LNCS, Springer, pp. 40–56. <https://eprint.iacr.org/2013/383/>. 1, 3
- [36] FLUHRER, S. Cryptanalysis of ring-LWE based key exchange with key share reuse. IACR Cryptology ePrint Archive report 2016/085, 2016. <http://eprint.iacr.org/2016/085>. 4

- [37] FUJIOKA, A., SUZUKI, K., XAGAWA, K., AND YONEYAMA, K. Practical and post-quantum authenticated key exchange from one-way secure key encapsulation mechanism. In *Symposium on Information, Computer and Communications Security, ASIA CCS 2013* (2013), K. Chen, Q. Xie, W. Qiu, N. Li, and W. Tzeng, Eds., ACM, pp. 83–94. [3](#)
- [38] GALBRAITH, S. D. Space-efficient variants of cryptosystems based on learning with errors, 2013. <https://www.math.auckland.ac.nz/~sgal018/compact-LWE.pdf>. [5](#)
- [39] GAMA, N., NGUYEN, P. Q., AND REGEV, O. Lattice enumeration using extreme pruning. In *Advances in Cryptology – EUROCRYPT 2010* (2010), H. Gilbert, Ed., vol. 6110 of *LNCS*, Springer, pp. 257–278. <http://www.iacr.org/archive/eurocrypt2010/66320257/66320257.pdf>. [8](#)
- [40] GARG, S., GENTRY, C., AND HALEVI, S. Candidate multilinear maps from ideal lattices. In *Advances in Cryptology – EUROCRYPT 2013* (2013), vol. 7881 of *LNCS*, Springer, pp. 1–17. <https://eprint.iacr.org/2012/610>. [8](#)
- [41] GENTLEMAN, W. M., AND SANDE, G. Fast Fourier transforms: for fun and profit. In *Fall Joint Computer Conference* (1966), vol. 29 of *AFIPS Proceedings*, pp. 563–578. http://cis.rit.edu/class/sing716/FFT_Fun_Profit.pdf. [11](#)
- [42] GENTRY, C. Fully homomorphic encryption using ideal lattices. In *STOC '09 Proceedings of the forty-first annual ACM symposium on Theory of computing* (2009), ACM, pp. 169–178. <https://www.cs.cmu.edu/~odonnell/hits09/gentry-homomorphic-encryption.pdf>. [1](#)
- [43] GENTRY, C., PEIKERT, C., AND VAIKUNTANATHAN, V. Trappeddoors for hard lattices and new cryptographic constructions. In *STOC '08 Proceedings of the fortieth annual ACM symposium on Theory of computing* (2008), ACM, pp. 197–206. <https://eprint.iacr.org/2007/432/>. [5](#)
- [44] GENTRY, C., AND SZYDLO, M. Cryptanalysis of the revised NTRU signature scheme. In *Advances in Cryptology – EUROCRYPT 2002* (2002), L. R. Knudsen, Ed., vol. 2332 of *LNCS*, Springer, pp. 299–320. https://www.iacr.org/archive/eurocrypt2002/23320295/nsssign_short3.pdf. [8](#)
- [45] GHOSH, S., AND KATE, A. Post-quantum secure onion routing (future anonymity in today’s budget). IACR Cryptology ePrint Archive report 2015/008, 2015. <http://eprint.iacr.org/2015/008>. [3](#)
- [46] GÖTTERT, N., FELLER, T., SCHNEIDER, M., BUCHMANN, J. A., AND HUSS, S. A. On the design of hardware building blocks for modern lattice-based encryption schemes. In *Cryptographic Hardware and Embedded Systems - CHES 2012* (2012), E. Prouff and P. Schaumont, Eds., vol. 7428 of *LNCS*, Springer, pp. 512–529. <http://www.iacr.org/archive/ches2012/74280511/74280511.pdf>. [7, 8, 9, 10](#)
- [47] GUERON, S. Parallelized hashing via j-lanes and j-pointers tree modes, with applications to SHA-256. IACR Cryptology ePrint Archive report 2014/170, 2014. <https://eprint.iacr.org/2014/170>. [12](#)
- [48] GUERON, S., AND KRASNOV, V. Fast prime field elliptic-curve cryptography with 256-bit primes. *Journal of Cryptographic Engineering* 5, 2 (2014), 141–151. <https://eprint.iacr.org/2013/816/>. [13](#)
- [49] GUERON, S., AND SCHLIEKER, F. Speeding up R-LWE post-quantum key exchange. IACR Cryptology ePrint Archive report 2016/467, 2016. <https://eprint.iacr.org/2016/467>. [10, 14](#)
- [50] GÜNEYSU, T., ODER, T., PÖPPELMANN, T., AND SCHWABE, P. Software speed records for lattice-based signatures. In *Post-Quantum Cryptography* (2013), P. Gaborit, Ed., vol. 7932 of *LNCS*, Springer, pp. 67–82. <http://cryptojedi.org/papers/#lattisigns>. [10, 13](#)
- [51] HANROT, G., PUJOL, X., AND STEHLÉE, D. Terminating BKZ. IACR Cryptology ePrint Archive report 2011/198, 2011. <http://eprint.iacr.org/2011/198/>. [8](#)
- [52] HOFFSTEIN, J., PIPHER, J., SCHANCK, J. M., SILVERMAN, J. H., AND WHYTE, W. Practical signatures from the partial Fourier recovery problem. In *Applied Cryptography and Network Security* (2014), I. Boureau, P. Owesarski, and S. Vaudey, Eds., vol. 8479 of *LNCS*, Springer, pp. 476–493. <https://eprint.iacr.org/2013/757/>. [1](#)
- [53] HOFFSTEIN, J., PIPHER, J., SCHANCK, J. M., SILVERMAN, J. H., WHYTE, W., AND ZHANG, Z. Choosing parameters for NTRUEncrypt. IACR Cryptology ePrint Archive report 2015/708, 2015. <http://eprint.iacr.org/2015/708>. [9](#)
- [54] HOFFSTEIN, J., PIPHER, J., AND SILVERMAN, J. H. NTRU: a ring-based public key cryptosystem. In *Algorithmic number theory* (1998), J. P. Buhler, Ed., vol. 1423 of *LNCS*, Springer, pp. 267–288. <https://www.securityinnovation.com/uploads/Crypto/ANTS97.ps.gz>. [1, 4](#)
- [55] HOWGRAVE-GRAHAM, N. A hybrid lattice-reduction and meet-in-the-middle attack against NTRU. In *Advances in Cryptology – CRYPTO 2007*, A. Menezes, Ed., vol. 4622 of *LNCS*. Springer, 2007, pp. 150–169. <http://www.iacr.org/archive/crypto2007/46220150/46220150.pdf>. [9](#)
- [56] KIRCHNER, P., AND FOUQUE, P.-A. An improved BKW algorithm for LWE with applications to cryptography and lattices. In *Advances in Cryptology – CRYPTO 2015* (2015), R. Gennaro and M. Robshaw, Eds., vol. 9215 of *LNCS*, Springer, pp. 43–62. <https://eprint.iacr.org/2015/552/>. [8, 18](#)
- [57] LAARHOVEN, T. *Search problems in cryptography*. PhD thesis, Eindhoven University of Technology, 2015. <http://www.thijs.com/docs/phd-final.pdf>. [8](#)
- [58] LAARHOVEN, T. Sieving for shortest vectors in lattices using angular locality-sensitive hashing. In *Advances in Cryptology – CRYPTO 2015* (2015), R. Gennaro and M. Robshaw, Eds., vol. 9216 of *LNCS*, Springer, pp. 3–22. <https://eprint.iacr.org/2014/744/>. [8](#)
- [59] LAARHOVEN, T., MOSCA, M., AND VAN DE POL, J. Finding shortest lattice vectors faster using quantum search. *Designs, Codes and Cryptography* 77, 2 (2015), 375–400. <https://eprint.iacr.org/2014/907/>. [8](#)
- [60] LANGLEY, A. TLS symmetric crypto. Blog post on imperialviolet.org, 2014. <https://www.imperialviolet.org/2014/02/27/tlssymmetriccrypto.html> (accessed 2015-10-07). [12](#)
- [61] LANGLEY, A., AND CHANG, W.-T. ChaCha20 and Poly1305 based cipher suites for TLS: Internet draft. <https://tools.ietf.org/html/draft-agl-tls-chacha20poly1305-04> (accessed 2015-02-01). [12](#)
- [62] LANGLEY, A., HAMBURG, M., AND TURNER, S. RFC 7748: Elliptic curves for security, 2016. <https://www.rfc-editor.org/rfc/rfc7748.txt>. [13](#)
- [63] LENSTRA, H. W., AND SILVERBERG, A. Revisiting the Gentry-Szydlo algorithm. In *Advances in Cryptology – CRYPTO 2014*, J. A. Garay and R. Gennaro, Eds., vol. 8616 of *LNCS*. Springer, 2014, pp. 280–296. <https://eprint.iacr.org/2014/430>. [8](#)
- [64] LEYS, J., GHYS, E., AND ALVAREZ, A. Dimensions, 2010. <http://www.dimensions-math.org/> (accessed 2015-10-19). [7, 18](#)
- [65] LINDNER, R., AND PEIKERT, C. Better key sizes (and attacks) for LWE-based encryption. In *Topics in Cryptology - CT-RSA 2011* (2011), A. Kiayias, Ed., vol. 6558 of *LNCS*, Springer, pp. 319–339. <https://eprint.iacr.org/2010/613/>. [3](#)

- [66] LIU, Z., SEO, H., ROY, S. S., GROSSCHÄDL, J., KIM, H., AND VERBAUWHEDE, I. Efficient Ring-LWE encryption on 8-bit AVR processors. In *Cryptographic Hardware and Embedded Systems - CHES 2015* (2015), T. Güneysu and H. Handschuh, Eds., vol. 9293 of *LNCS*, Springer, pp. 663–682. <https://eprint.iacr.org/2015/410/>. 3, 7, 10
- [67] LONGA, P., AND NAEHRIG, M. Speeding up the number theoretic transform for faster ideal lattice-based cryptography. IACR Cryptology ePrint Archive report 2016/504, 2016. <https://eprint.iacr.org/2016/504>. 14
- [68] LYUBASHEVSKY, V. Lattice signatures without trapdoors. In *Advances in Cryptology – EUROCRYPT 2012* (2012), D. Pointcheval and T. Johansson, Eds., vol. 7237 of *LNCS*, Springer, pp. 738–755. <https://eprint.iacr.org/2011/537/>. 5
- [69] LYUBASHEVSKY, V., PEIKERT, C., AND REGEV, O. On ideal lattices and learning with errors over rings. In *Advances in Cryptology – EUROCRYPT 2010* (2010), H. Gilbert, Ed., vol. 6110 of *LNCS*, Springer, pp. 1–23. <http://www.di.ens.fr/~lyubash/papers/ringLWE.pdf>. 2
- [70] LYUBASHEVSKY, V., PEIKERT, C., AND REGEV, O. On ideal lattices and learning with errors over rings. *Journal of the ACM (JACM)* 60, 6 (2013), 43:1–43:35. <http://www.cims.nyu.edu/~regev/papers/ideal-lwe.pdf>. 1, 5
- [71] MELCHOR, C. A., BARRIER, J., FOUSSE, L., AND KILLIJIAN, M. XPIRe: Private information retrieval for everyone. IACR Cryptology ePrint Archive report 2014/1025, 2014. <http://eprint.iacr.org/2014/1025>. 10
- [72] MICCIANCIO, D., AND VOULGARIS, P. Faster exponential time algorithms for the shortest vector problem. In *SODA '10 Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms* (2010), SIAM, pp. 1468–1480. <http://dl.acm.org/citation.cfm?id=1873720>. 8
- [73] MONTANARO, A. Quantum walk speedup of backtracking algorithms. arXiv preprint arXiv:1509.02374, 2015. <http://arxiv.org/pdf/1509.02374v2>. 8
- [74] MONTGOMERY, P. L. Modular multiplication without trial division. *Mathematics of Computation* 44, 170 (1985), 519–521. <http://www.ams.org/journals/mcom/1985-44-170/S0025-5718-1985-0777282-X/S0025-5718-1985-0777282-X.pdf>. 11
- [75] NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY. FIPS PUB 202 – SHA-3 standard: Permutation-based hash and extendable-output functions, 2015. <http://nvlpubs.nist.gov/nistpubs/FIPS/NIST.FIPS.202.pdf>. 5, 10
- [76] NATIONAL INSTITUTE OF STANDARDS AND TECHNOLOGY. Workshop on cybersecurity in a post-quantum world, 2015. <http://www.nist.gov/itl/csd/ct/post-quantum-crypto-workshop-2015.cfm>. 1
- [77] NATIONAL SECURITY AGENCY. NSA suite B cryptography. https://www.nsa.gov/ia/programs/suiteb_cryptography/, Updated on August 19, 2015. 1
- [78] NGUYEN, P. Q., AND VIDICK, T. Sieve algorithms for the shortest vector problem are practical. *Journal of Mathematical Cryptology* 2, 2 (2008), 181–207. <ftp://ftp.di.ens.fr/pub/users/pnguyen/JoMC08.pdf>. 8
- [79] NUSSBAUMER, H. J. Fast polynomial transform algorithms for digital convolution. *IEEE Transactions on Acoustics, Speech and Signal Processing* 28, 2 (1980), 205–215. 13
- [80] PEIKERT, C. Lattice cryptography for the Internet. In *Post-Quantum Cryptography* (2014), M. Mosca, Ed., vol. 8772 of *LNCS*, Springer, pp. 197–219. <http://web.eecs.umich.edu/~cpeikert/pubs/suiteb.pdf>. 1, 2, 3, 4, 19, 21
- [81] PÖPPELMANN, T., DUCAS, L., AND GÜNEYSU, T. Enhanced lattice-based signatures on reconfigurable hardware. In *Cryptographic Hardware and Embedded Systems – CHES 2014* (2014), L. Batina and M. Robshaw, Eds., vol. 8731 of *LNCS*, Springer, pp. 353–370. <https://eprint.iacr.org/2014/254/>. 3
- [82] PÖPPELMANN, T., AND GÜNEYSU, T. Towards practical lattice-based public-key encryption on reconfigurable hardware. In *Selected Areas in Cryptography – SAC 2013* (2013), T. Lange, K. Lauter, and P. Lisoněk, Eds., vol. 8282 of *LNCS*, Springer, pp. 68–85. https://www.ei.rub.de/media/sh/veroeffentlichungen/2013/08/14/lwe_encrypt.pdf. 1, 6, 7, 10, 11
- [83] REGEV, O. On lattices, learning with errors, random linear codes, and cryptography. *Journal of the ACM* 56, 6 (2009), 34. <http://www.cims.nyu.edu/~regev/papers/qcrypto.pdf>. 5
- [84] RÉNYI, A. On measures of entropy and information. In *Fourth Berkeley symposium on mathematical statistics and probability* (1961), vol. 1, pp. 547–561. <http://projecteuclid.org/euclid.bsm/1200512181>. 18
- [85] ROY, S. S., VERCAUTEREN, F., MENTENS, N., CHEN, D. D., AND VERBAUWHEDE, I. Compact Ring-LWE cryptoprocessor. In *Cryptographic Hardware and Embedded Systems – CHES 2014* (2014), L. Batina and M. Robshaw, Eds., vol. 8731 of *LNCS*, Springer, pp. 371–391. <https://eprint.iacr.org/2013/866/>. 3, 5, 11
- [86] SCHNORR, C.-P., AND EUCHNER, M. Lattice basis reduction: improved practical algorithms and solving subset sum problems. *Mathematical programming* 66, 1-3 (1994), 181–199. http://www.csie.nuk.edu.tw/~cychen/Lattices/Lattice%20Basis%20Reduction_%20Improved%20Practical%20Algorithms%20and%20Solving%20Subset%20Sum%20Problems.pdf. 8
- [87] STEHLÉ, D., AND STEINFELD, R. Making NTRU as secure as worst-case problems over ideal lattices. In *Advances in Cryptology – EUROCRYPT 2011* (2011), K. G. Paterson, Ed., vol. 6632 of *LNCS*, Springer, pp. 27–47. <http://www.iacr.org/archive/eurocrypt2011/66320027/66320027.pdf>. 1, 4
- [88] Tor project: Anonymity online. <https://www.torproject.org/>. 3
- [89] VERSHYNIN, R. Introduction to the non-asymptotic analysis of random matrices. In *Compressed sensing*. Cambridge University Press, 2012, pp. 210–268. <http://www-personal.umich.edu/~romanv/papers/non-asymptotic-rmt-plain.pdf>. 20
- [90] ZHANG, J., ZHANG, Z., DING, J., SNOOK, M., AND DAGDELEN, Ö. Authenticated key exchange from ideal lattices. In *Advances in Cryptology – EUROCRYPT 2015* (2015), E. Oswald and M. Fischlin, Eds., vol. 9057 of *LNCS*, Springer, pp. 719–751. <https://eprint.iacr.org/2014/589/>. 3

A JARJAR

As discussed in Section 6, it is also possible to instantiate our scheme with dimension $n = 512$, modulus $q = 12289$ and noise parameter $k = 24$, to obtain a decent security level against currently known attacks: our analysis leads to a claim of 118-bits of post quantum security. Most of the scheme is similar, except that one should use \tilde{D}_2 as reconciliation lattice, and the provable failure bound drops to 2^{-55} . It is quite likely that a more aggressive

analysis would lead to a claim of 128 bits of security against the best known attacks.

The performance of almost all computations in NEWHOPE scales linearly in n (except for the NTT, which scales quasi-linear and a small constant overhead, for example, for reading a seed from `/dev/urandom`). Also, noise sampling scales linearly in n and k , so JARJAR with $k = 24$ is expected to take $3/4$ of the time for sampling of a noise polynomial compared to NEWHOPE. Message sizes scale linearly in n (except for the constant-size 32-byte seed). The performance of JARJAR is thus expected to be about a factor of 2 better in terms of size and slightly less than a factor of 2 in terms of speed. We confirmed this by adapting our implementations of NEWHOPE to the JARJAR parameters. The C reference implementation of JARJAR takes 167 102 Haswell cycles for the key generation on the server, 234 268 Haswell cycles for the key generation and joint-key computation on the client, and 45 712 Haswell cycles for the joint-key computation on the server. The AVX2 optimized implementation of JARJAR takes 71 608 Haswell cycles for the key generation on the server, 84 316 Haswell cycles for the key generation and joint-key computation on the client, and 12 944 Haswell cycles for the joint-key computation on the server. As stated before, we do not recommend to use JARJAR, but only provide these numbers as a target for comparison to other more aggressive RLWE cryptosystems in the literature.

B Proof of Theorem 4.1

A simple security reduction to rounded Gaussians. In [9], Bai et al. identify Rényi divergence as a powerful tool to improve or generalize security reductions in lattice-based cryptography. We review the key properties. The Rényi divergence [9, 84] is parametrized by a real $a > 1$, and defined for two distributions P, Q by:

$$R_a(P||Q) = \left(\sum_{x \in \text{Supp}(P)} \frac{P(x)^a}{Q(x)^{a-1}} \right)^{\frac{1}{a-1}}.$$

It is multiplicative: if P, P' are independents, and Q, Q' are also independents, then $R_a(P \times P' || Q \times Q') \leq R_a(P||Q) \cdot R_a(P'||Q')$. Finally, Rényi divergence relates the probabilities of the same event E under two different distributions P and Q :

$$Q(E) \geq P(E)^{a/(a-1)} / R_a(P||Q).$$

For our argument, recall that because the final shared key μ is obtained through hashing as $\mu \leftarrow \text{SHA3-256}(v)$ before being used, then, in the random oracle model (ROM), any successful attacker must recover v exactly.

We call this event E . We also define ξ to be the rounded Gaussian distribution of parameter $\sigma = \sqrt{k/2} = \sqrt{8}$, that is the distribution of $\lfloor \sqrt{8} \cdot x \rfloor$ where x follows the standard normal distribution.

A simple script (`scripts/Renyi.py`) computes $R_9(\psi_{16}||\xi) \approx 1.00063$. Yet because $5n = 5120$ samples are used per instance of the protocol, we need to consider the divergence $R_9(P||Q) = R_9(\psi_{16}, \xi)^{5n} \leq 26$ where $P = \psi_{16}^{5n}$ and $Q = \xi^{5n}$. We conclude as follows.

The choice $a = 9$ is rather arbitrary but seemed a good trade-off between the coefficient $1/R_a(\psi_{16}||\xi)$ and the exponent $a/(a-1)$. This reduction is provided as a safeguard: switching from Gaussian to binomial distributions can not dramatically decrease the security of the scheme. With practicality in mind, we will simply ignore the loss factor induced by the above reduction, since the best-known attacks against LWE do not exploit the structure of the error distribution, and seem to depend only on the standard deviation of the error (except in extreme cases [6, 56]).

C Details on Reconciliation

In this section, we provide more technical details on the recovery mechanism involving the lattice \tilde{D}_4 . For reader's convenience, some details are repeated from Section 5.

Splitting for recovery. By $\mathcal{S} = \mathbb{Z}[X]/(X^4 + 1)$ we denote the 8-th cyclotomic ring, having rank 4 over \mathbb{Z} . Recall that our full ring is $\mathcal{R} = \mathbb{Z}[X]/(X^n + 1)$ for $n = 1024$. An element of the full ring \mathcal{R} can be identified to a vector $(f'_0, \dots, f'_{255}) \in \mathcal{S}^{n/4}$ such that

$$f(X) = f'_0(X^{256}) + X f'_1(X^{256}) + \dots + X^{255} f'_{255}(X^{256}).$$

In other words, the coefficients of f'_i are the coefficients $f_i, f_{i+256}, f_{i+512}, f_{i+768}$.

The lattice \tilde{D}_4 . One may construct the lattice \tilde{D}_4 as two shifted copies of \mathbb{Z}^4 using a glue vector \mathbf{g} :

$$\tilde{D}_4 = \mathbb{Z}^4 \cup \mathbf{g} + \mathbb{Z}^4 \text{ where } \mathbf{g}' = \left(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2} \right).$$

We recall that \tilde{D}_4 provides the densest lattice sphere packing in dimension 4 [30], this suggest it is optimal for error correction purposes. The Voronoi cell \mathcal{V} of \tilde{D}_4 is the icositetrachoron [64] (a.k.a. the 24-cell, the convex regular 4-polytope with 24 octahedral cells) and the Voronoi relevant vectors of two types: 8 type-A vectors $(\pm 1, 0, 0, 0), (0, \pm 1, 0, 0), (0, 0, \pm 1, 0), (0, 0, 0, \pm 1)$, and 16 type-B vectors $(\pm \frac{1}{2}, \pm \frac{1}{2}, \pm \frac{1}{2}, \pm \frac{1}{2})$. The natural condition to correct decoding in \tilde{D}_4 should therefore be $e \in \mathcal{V}$, and this can be expressed as $\langle e, v \rangle \leq$

1/2 for all Voronoi relevant vectors v . Interestingly, those 24-linear inequalities can be split as $\|e\|_1 \leq 1$ (providing the 16 inequalities for the type-B vectors) and $\|e\|_\infty \leq 1/2$ (providing the 8 inequalities for the type-A vectors). In other words, the 24-cell \mathcal{V} is the intersection of an ℓ_1 -ball (an hexadecachoron) and an ℓ_∞ -ball (a tesseract).

As our basis for \tilde{D}_4 , we will choose $\mathbf{B} = (\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2, \mathbf{g})$ where \mathbf{u}_i are the canonical basis vectors of \mathbb{Z}^4 . The construction of \tilde{D}_4 with a glue vector gives a simple and efficient algorithm for finding closest vectors in \tilde{D}_4 . Note that we have $\mathbf{u}_3 = -\mathbf{u}_0 - \mathbf{u}_1 - \mathbf{u}_2 + 2\mathbf{g} = \mathbf{B} \cdot (-1, -1, -1, 2)^t$.

Algorithm 1 (restated) $\text{CVP}_{\tilde{D}_4}(\mathbf{x} \in \mathbb{R}^4)$

Ensure: An integer vector \mathbf{z} such that \mathbf{Bz} is a closest vector to \mathbf{x} : $\mathbf{x} - \mathbf{Bz} \in \mathcal{V}$

- 1: $\mathbf{v}_0 \leftarrow \lfloor \mathbf{x} \rfloor$
 - 2: $\mathbf{v}_1 \leftarrow \lfloor \mathbf{x} - \mathbf{g} \rfloor$
 - 3: $k \leftarrow (\|\mathbf{x} - \mathbf{v}_0\|_1 < 1) ? 0 : 1$
 - 4: $(v_0, v_1, v_2, v_3)^t \leftarrow \mathbf{v}_k$
 - 5: **return** $(v_0, v_1, v_2, k)^t + v_3 \cdot (-1, -1, -1, 2)^t$
-

Decoding in \tilde{D}_4/\mathbb{Z}^4 . Because $\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2$ and $2\mathbf{g}$ belong to \mathbb{Z}^4 , a vector in \tilde{D}_4/\mathbb{Z}^4 is simply given by the parity of its last coordinate in base \mathbf{B} . This gives an even simpler algorithm to encode and decode a bit in \tilde{D}_4/\mathbb{Z}^4 . The encoding is given by $\text{Encode}(k \in \{0, 1\}) = k\mathbf{g}$ and the decoding is given below.

Algorithm 2 (restated) $\text{Decode}(\mathbf{x} \in \mathbb{R}^4/\mathbb{Z}^4)$

Ensure: A bit k such that $k\mathbf{g}$ is a closest vector to $\mathbf{x} + \mathbb{Z}^4$: $\mathbf{x} - k\mathbf{g} \in \mathcal{V} + \mathbb{Z}^4$

- 1: $\mathbf{v} = \mathbf{x} - \lfloor \mathbf{x} \rfloor$
 - 2: **return** 0 if $\|\mathbf{v}\|_1 \leq 1$ and 1 otherwise
-

When we want to decode to \tilde{D}_4/\mathbb{Z}^4 rather than to \tilde{D}_4 , the 8 inequalities given by type-A vectors are irrelevant since those vectors belong to \mathbb{Z}^4 . It follows that:

Lemma C.1 *For any $k \in \{0, 1\}$ and any $\mathbf{e} \in \mathbb{R}^4$ such that $\|\mathbf{e}\|_1 < 1$, we have $\text{Decode}(k\mathbf{g} + \mathbf{e}) = k$.*

C.1 Reconciliation

We define the following r -bit reconciliation function:

$$\text{HelpRec}(\mathbf{x}; b) = \text{CVP}_{\tilde{D}_4} \left(\frac{2^r}{q} (\mathbf{x} + b\mathbf{g}) \right) \bmod 2^r,$$

where $b \in \{0, 1\}$ is a uniformly chosen random bit. This random vector is equivalent to the ‘‘doubling’’ trick of

Peikert [80]. Indeed, because q is odd, it is not possible to map deterministically the uniform distribution from \mathbb{Z}_q^4 to \mathbb{Z}_2 , which necessary results in a final bias.

Lemma C.2 *Assume $r \geq 1$ and $q \geq 9$. For any $\mathbf{x} \in \mathbb{Z}_q^4$, set $\mathbf{r} := \text{HelpRec}(\mathbf{x}) \in \mathbb{Z}_{2^r}^4$. Then, $\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br} \bmod 1$ is close to a point of \tilde{D}_4/\mathbb{Z}^4 , precisely, for $\alpha = 1/2^r + 2/q$:*

$$\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br} \in \alpha\mathcal{V} + \mathbb{Z}^4 \quad \text{or} \quad \frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br} \in \mathbf{g} + \alpha\mathcal{V} + \mathbb{Z}^4.$$

Additionally, for \mathbf{x} uniformly chosen in \mathbb{Z}_q^4 we have $\text{Decode}(\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br})$ is uniform in $\{0, 1\}$ and independent of \mathbf{r} .

Proof One can easily check the correctness of the $\text{CVP}_{\tilde{D}_4}$ algorithm: for any \mathbf{y} , $\mathbf{y} - \mathbf{B} \cdot \text{CVP}_{\tilde{D}_4}(\mathbf{y}) \in \mathcal{V}$. To conclude with the first property, it remains to note that $\mathbf{g} \in 2\mathcal{V}$, and that \mathcal{V} is convex.

For the second property, we show that there is a permutation $\pi : (\mathbf{x}, b) \mapsto (\mathbf{x}', b')$ of $\mathbb{Z}_q^4 \times \mathbb{Z}_2$, such that, for all (\mathbf{x}, b) it holds that:

$$\text{HelpRec}(\mathbf{x}; b) = \text{HelpRec}(\mathbf{x}'; b') \quad (= \mathbf{r}) \quad (3)$$

$$\text{Decode} \left(\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br} \right) = \text{Decode} \left(\frac{1}{q}\mathbf{x}' - \frac{1}{2^r}\mathbf{Br} \right) \oplus 1 \quad (4)$$

where \oplus denotes the xor operation. We construct π as follows: set $b' := b \oplus 1$ and $\mathbf{x}' = \mathbf{x} + (b - b' + q)\mathbf{g} \bmod q$. Note that $b - b' + q$ is always even so $(b - b' + q)\mathbf{g}$ is always well defined in \mathbb{Z}_q (recall that $\mathbf{g} = (1/2, 1/2, 1/2, 1/2)^t$). It follows straightforwardly that $\frac{2^r}{q}(\mathbf{x} + b\mathbf{g}) - \frac{2^r}{q}(\mathbf{x}' + b'\mathbf{g}) = -2^r\mathbf{g} \bmod 2^r$. Since $\mathbf{g} \in \tilde{D}_4$, condition (3) holds. For condition (4), notice that:

$$\begin{aligned} \frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br} &\in k\mathbf{g} + \alpha\mathcal{V} \quad \bmod 1 \\ \implies \frac{1}{q}\mathbf{x}' - \frac{1}{2^r}\mathbf{Br} &\in (k \oplus 1)\mathbf{g} + \alpha'\mathcal{V} \quad \bmod 1 \end{aligned}$$

for $\alpha' = \alpha + 2/q$. Because $r \geq 1$ and $q \geq 9$, we have $\alpha' = 1/2^r + 4/q < 1$, and remembering that $\mathbf{e} \in \mathcal{V} \implies \|\mathbf{e}\|_1 \leq 1$ (inequalities for type-B vectors), one concludes by Lemma C.1. ■

It remains to define $\text{Rec}(\mathbf{x}, \mathbf{r}) = \text{Decode}(\frac{1}{q}\mathbf{x} - \frac{1}{2^r}\mathbf{Br})$ to describe a 1-bit-out-of-4-dimensions reconciliation protocol (Protocol 4). Those functions are extended to 256-bits out of 1024-dimensions by the splitting described at the beginning of this section.

Alice		Bob
$\mathbf{x}' \in \mathbb{Z}_q^4$	$\mathbf{x}' \approx \mathbf{x}$	$\mathbf{x} \in \mathbb{Z}_q^4$
	$\xleftarrow{\mathbf{r}}$	$\mathbf{r} \leftarrow \text{HelpRec}(\mathbf{x}) \in \mathbb{Z}_{2^r}^4$
$k' \leftarrow \text{Rec}(\mathbf{x}', \mathbf{r})$		$k \leftarrow \text{Rec}(\mathbf{x}, \mathbf{r})$

Protocol 4: Reconciliation protocol in $q\tilde{D}_4/q\mathbb{Z}^4$.

Lemma C.3 *If $\|\mathbf{x} - \mathbf{x}'\|_1 < (1 - 1/2^r) \cdot q - 2$, then by the above protocol $4k = k'$. Additionally, if \mathbf{x} is uniform, then k is uniform independently of \mathbf{r} .*

Fixed-point implementation. One remarks that, while we described our algorithm in \mathbb{R} for readability, floating-point-arithmetic is not required in practice. Indeed, all computation can be performed using integer arithmetic modulo $2^r q$. Our parameters ($r = 2, q = 12289$) are such that $2^r q < 2^{16}$, which offers a good setting also for small embedded devices.

D Analysis of the failure probability

To proceed with the task of bounding—as tightly as possible—the failure probability, we rely on the notion of moments of a distribution and of subgaussian random variables [89]. We recall that the moment-generating function of a real random variable \mathcal{X} is defined as follows:

$$M_{\mathcal{X}}(t) := \mathbb{E}[\exp(t(\mathcal{X} - \mathbb{E}[\mathcal{X}]))].$$

We extend the definition to distributions over \mathbb{R} : $M_{\phi} := M_{\mathcal{X}}$ where $\mathcal{X} \leftarrow \phi$. Note that $M_{\phi}(t)$ is not necessary finite for all t , but it is the case if the support of ϕ is bounded. We also recall that if \mathcal{X} and \mathcal{Y} are independent, then the moment-generating functions verify the identity $M_{\mathcal{X}+\mathcal{Y}}(t) = M_{\mathcal{X}}(t) \cdot M_{\mathcal{Y}}(t)$.

Theorem D.1 (Chernoff-Cramer inequality) *Let ϕ be a distribution over \mathbb{R} and let $\mathcal{X}_1 \dots \mathcal{X}_n$ be i.i.d. random variables of law ϕ , with average μ . Then, for any t such that $M_{\phi}(t) < \infty$ it holds that*

$$\mathbb{P}\left[\sum_{i=1}^n \mathcal{X}_i \geq n\mu + \beta\right] \leq \exp(-\beta t + n \ln(M_{\phi}(t))).$$

Definition A centered distribution ϕ over \mathbb{R} is said to be σ -subgaussian if its moment-generating function verifies $\mathbb{E}_{\mathcal{X} \leftarrow \phi}[\exp(t\mathcal{X})] \leq \exp(2t^2\sigma^2)$.

A special case of Chernoff-Cramer bound follows by choosing t appropriately.

Lemma D.2 (Adapted from [89]) *If \mathbf{x} has independently chosen coordinates from ϕ , a σ -subgaussian distribution, then, for any vector $\mathbf{v} \in \mathbb{R}^n$, except with probability less than $\exp(-\tau^2/2)$, we have:*

$$\langle \mathbf{x}, \mathbf{v} \rangle \leq \|\mathbf{v}\| \sigma \tau.$$

The centered binomial distribution ψ_k of parameter k is $\sqrt{k/2}$ -subgaussian. This is established from the fact that $b_0 - b'_0$ is $\frac{1}{2}$ -subgaussian (which is easy to check),

and by Euclidean additivity of k independent subgaussian variables. Therefore, the binomial distribution used (of parameter $k = 16$) is $\sqrt{8}$ -subgaussian.

We here propose a rather tight tail-bound on the error term. We recall that the difference \mathbf{d} in the agreed key before key reconciliation is $\mathbf{d} = \mathbf{e}\mathbf{s}' - \mathbf{e}'\mathbf{s} + \mathbf{e}''$. We wish to bound $\|d'_i\|_1$ for all $i \leq 255$, where the $d'_i \in \mathcal{S}$ form the decomposition of d described in Section 5. Note that, seen as a vector over \mathbb{R} , we have

$$\|x\|_1 = \max_y \langle x, y \rangle,$$

where y ranges over $\{\pm 1\}^4$. We also remark that for $\mathbf{a}, \mathbf{b} \in \mathcal{R}$, one may rewrite $(\mathbf{ab})_i \in \mathcal{S}$

$$(\mathbf{ab})_i = \sum_{j=0}^{255} \pm a_j b_{(i-j) \bmod 256},$$

where the sign \pm depends only on the indices i, j . This allows to rewrite

$$\|(\mathbf{e}\mathbf{s}' - \mathbf{e}'\mathbf{s})_i\|_1 = \max_y \sum_{j=0}^{255} \pm (\langle e_i, s_{i-j}y \rangle + \langle e'_i, s'_{i-j}y \rangle) \quad (5)$$

where $y \in \mathcal{S}$ ranges over all polynomials with ± 1 coefficients.

Lemma D.3 *Let $s, s' \in \mathcal{S}$ be drawn with independent coefficients from ψ_{16} , then, except with probability 2^{-64} , the vectors $\mathbf{v}_y = (s_0y, \dots, s_{255}y, s'_0y, \dots, s'_{255}y) \in \mathbb{Z}^{2048}$ verify simultaneously $\|\mathbf{v}_y\|_2^2 \leq 102500$ for all $y \in \mathcal{S}$ with ± 1 coefficients.*

Proof Let us first remark that $\|\mathbf{v}\|_2^2 = \sum_{i=0}^{512} \|s_i y\|_2^2$ where the s_i 's are i.i.d. random variables following distribution ψ_{16}^4 . Because the support of ψ_{16}^4 is reasonably small (of size $(216+1)^4 \approx 2^{20}$), we numerically compute the distribution $\phi_y : \|sy\|_2^2$ where $s \leftarrow \psi_{16}^4$ for each y (see `scripts/failure.py`). Note that such numerical computation of the probability density function does not raise numerical stability concern, because it involves a depth-2 circuit of multiplication and addition of *positive* reals: the relative error growth remain at most quadratic in the length of the computation.

From there, one may compute $\mu = \mathbb{E}_{\mathcal{X} \leftarrow \phi_y}[\mathcal{X}]$ and $M_{\phi_y}(t)$ (similarly this computation has polynomial relative error growth assuming \exp is computed with constant relative error growth).

We apply Chernoff-Cramer inequality with parameters $n = 512, n\mu + \beta = 102500$ and $t = 0.0055$, and obtain the desired inequality for each given y , except with probability at most $2^{-68.56}$. We conclude by union-bound over the 2^4 choices of y . ■

Corollary D.4 For $\mathbf{s}, \mathbf{e}, \mathbf{e}', \mathbf{s}, \mathbf{s}' \in \mathcal{R}$ drawn with independent coefficients according to ψ_{16} , except with probability at most 2^{-61} we have simultaneously for all $i \leq 255$ that

$$\|(\mathbf{e}'\mathbf{s} + \mathbf{e}\mathbf{s}' + \mathbf{e}'')_i\|_1 \leq 9210 < \lfloor 3q/4 \rfloor - 2 = 9214.$$

Proof First, we know that $\|(\mathbf{e}'')_i\|_1 \leq 4 \cdot 16 = 64$ for all i 's, because the support of ψ is $[-16, 16]$. Now, for each i , write

$$\|(\mathbf{e}'\mathbf{s} + \mathbf{e}\mathbf{s}')_i\|_1 = \max_y \langle \mathbf{v}_{i,y}, \mathbf{e}^+ \rangle$$

according to Equation 5, where \mathbf{v}_y depends on \mathbf{s}, \mathbf{s}' , and \mathbf{e}^+ is a permutation of the 2048 coefficients of \mathbf{e} and \mathbf{e}' , each of them drawn independently from ψ_{16} . By Lemma D.3 $\|\mathbf{v}_{0,y}\|^2 \leq 102500$ for all y with probability less than 2^{-64} . Because $\mathbf{v}_{i,y}$ is equal to $\mathbf{v}_{0,y}$ up to a signed permutation of the coefficient, we have $\|\mathbf{v}_{i,y}\| \leq 102500$ for all i, y .

Now, for each i, y , we apply the tail bound of Lemma D.2 with $\tau = 10.1$ knowing that ψ_{16} is σ -subgaussian for $\sigma = \sqrt{16/2}$, and we obtain, except with probability at most $2 \cdot 2^{-73.5}$ that

$$|\langle \mathbf{v}_{i,y}, \mathbf{e} \rangle| \leq 9146.$$

By union bound, we conclude that the claimed inequality holds except with probability less than $2^{-64} + 2 \cdot 8 \cdot 256 \cdot 2^{-73} \leq 2^{-61}$. ■

E A small example of edge blurring

We illustrate what it means to “blur the edges” in the reconciliation mechanism (cf. Section 5). This is an adaptation of the 1-dimensional randomized doubling trick of Peikert [80]. Formally, the proof that this process results in an unbiased key agreement is provided in Lemma C.2.

Consider the two-dimensional example with $q = 9$ depicted in Figure 3. All possible values of a vector are depicted with red or black dots. All red dots (in the grey Voronoi cell) are mapped to 1 in the key bit; the black dots are mapped to a zero. There is a total of $9 \cdot 9 = 81$ dots, 40 red ones, and 41 black ones. This obviously leads to a biased key. To address this we use one random bit to determine whether we add the vector $(\frac{1}{2q}, \frac{1}{2q})$ before running error reconciliation. For most dots this makes no difference; they stay in the same Voronoi cell and the “directed noise” is negligible for larger values of q . However, some dots right at the border are moved to another Voronoi cell with probability $\frac{1}{2}$ (depending on the value of the random bit). In the example in Figure 3 we see that 72 dots (36 red and 36 black ones) remain in their Voronoi cell; the other 9 dots (4 red and 5 black ones) change Voronoi cell with probability $\frac{1}{2}$ which precisely eliminates the bias in the key.

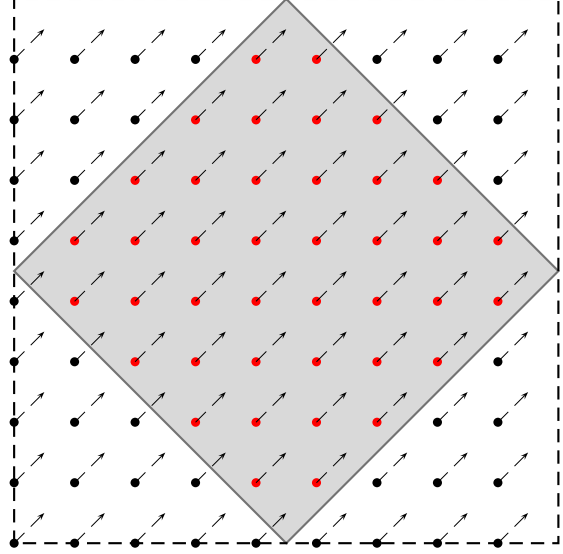


Figure 3: Possible values of a vector in $\mathbb{Z}_9 \times \mathbb{Z}_9$, their mapping to a zero or one bit, and the effect of blurring.

F Constant-time sampling of \mathbf{a}

For most applications, variable-time sampling of the public parameter \mathbf{a} is not a problem. This is simply because timing only leaks information that is part of the public messages. However, in anonymity networks like Tor, the situation is slightly different. In Tor, as part of the circuit establishment, a client performs key exchanges with three nodes. The first of these three nodes (the entry or guard node) knows the identity of the client, but already the second node does not; all communication, including the key exchange, with the second node is performed through the guard node. If the attacker has control over the second (or third) node of a Tor circuit *and* can run a timing attack on the client machine, then the attacker can check whether the timing information from the sampling of \mathbf{a} “matches” the public parameter received from that client and thus gain information about whether this is the client currently establishing a circuit. In other words, learning timing information does not break security of the key exchange, but it may break (or at least reduce) anonymity of the client.

For such scenarios we propose an efficient constant-time method to generate the public parameter \mathbf{a} . This constant-time approach is incompatible with the straightforward approach described in Section 7; applications have to decide what approach to use and then deploy it on all servers and clients. The constant-time approach works as follows:

1. Sample 16 blocks (2688 bytes) of SHAKE-128 output from the 32-byte seed. Consider this output an

84×16 -matrix of 16-bit unsigned little-endian integers $a_{i,j}$.

2. In each of 16 the columns of this matrix, move all entries that are $> 5q$ to the end without changing the order of the other elements (see below).
3. Check whether $a_{63,j}$ is smaller than $5q$ for $j = 0, \dots, 15$. In the extremely unlikely case that this does not hold, “squeeze” another 16 blocks of output from SHAKE-128 and go back to step 2.
4. Write the entries of the matrix as coefficients to the polynomial \mathbf{a} , where $a_{0,0}$ is the constant term, $a_{0,1}$ is the coefficient of X , $a_{0,2}$ is the coefficient of X^2 and so on, $a_{1,0}$ is the coefficient of X^{15} and so on and finally $a_{63,15}$ is the coefficient of X^{1023} .

The interesting part of this algorithm for constant-time implementations is the “move entries that are $> 5q$ to the end” step. This can efficiently be achieved by using a sorting network like Batchersort [11], where the usual compare-and-swap operator is replaced by an operator that swaps two elements if the one at the lower position is $\geq 5q$. As the columns are “sorted” independently, the sorting network can be computed 16 times in parallel, which is particularly efficient using the AVX2 vector instructions. We implemented this approach in AVX2 and noticed only a very small slowdown compared to the simpler approach described in Section 7: the median of 10000000 times sampling \mathbf{a} is 37292 cycles (average: 37310) on the same Haswell CPU that we used for all other benchmarks. For the reference implementation the slowdown is more noticeable: the median of 10000000 times sampling \mathbf{a} is 96552 cycles (average: 96562).