



McBits revisited: toward a fast constant-time code-based KEM

Tung Chou¹

Received: 19 December 2017 / Accepted: 2 March 2018 / Published online: 12 March 2018
© Springer-Verlag GmbH Germany, part of Springer Nature 2018

Abstract

This paper presents a constant-time fast implementation for a high-security code-based key encapsulation mechanism (KEM). The implementation is based on the “McBits” paper by Bernstein, Chou, and Schwabe in 2013: we use the same FFT algorithms for root finding and syndrome computation, similar algorithms for secret permutation, and bitslicing for low-level operations. As opposed to McBits, where a high decryption throughput is achieved by running many decryption operations in parallel, we take a different approach to exploit the internal parallelism in one decryption operation for the use of more applications. As a result, we manage to achieve a slightly better decryption throughput at a much higher security level than McBits. As a minor contribution, we also present a constant-time implementation for encryption and key-pair generation, with similar techniques used for decryption.

Keywords McEliece · Niederreiter · Bitslicing · Software implementation

1 Introduction

In recent years, due to the advance in quantum computing, cryptographers are paying more and more attention to post-quantum cryptography. In particular, NIST’s call for proposals [23] serves as an announcement to declare that post-quantum cryptography is going to be reality, and the whole world needs to be prepared for that. Among other things, we need post-quantum public-key encryption schemes and key encapsulation mechanisms (KEMs), and some of the most promising candidates today are from code-based cryptography.

In 1978, McEliece proposed his hidden-Goppa-code cryptosystem [20] as the first code-based encryption system. Until today, almost 40 years of research has been invested on cryptanalyzing the system, yet nothing has really shaken its security. It has thus become one of the most confidence-inspiring post-quantum encryption systems we have today, and it is important to evaluate how practical the system is for deployment.

In 2013, Bernstein, Chou, and Schwabe published the “McBits” paper [6], which presents a software implementation of Niederreiter’s dual form [22] of the McEliece

cryptosystem. McBits features (1) a very high decoding (and thus decryption) throughput which is an order of magnitude faster than the previous implementation by Biswas and Sendrier [12], and (2) full protection against timing attacks. These features are achieved by bitslicing non-conventional algorithms for decoding: they use the Gao–Mateer additive FFT [17] for the root finding, the corresponding “transposed” FFT for syndrome computation, and a sorting network for secret permutation.

The decryption throughput McBits achieves, however, relies on the assumption that there are many decryption operations that can be carried out at the same time. This is a reasonable assumption for some applications, but not for the all applications. The user would be glad to have an implementation that is capable of decrypting efficiently, even when there is only one decryption operation at the moment.

The main contribution of this paper is that we show the assumption is NOT a requirement to achieve a high decryption throughput. Even better, our software actually achieves a slightly better decryption throughput than McBits, at a much higher security level. To achieve this, we need to have a deep understanding about the data flow in each stage of decoding algorithm in order to figure out what kind of internal parallelism there is and how it can be exploited.

Note that this paper is an extended version of [14]. The implementation presented in this paper is mostly the same as that for [14]. The main difference lies in generation the

✉ Tung Chou
blueprint@crypto.tw

¹ Graduate School of Engineering, Osaka University, 1-1 Yamadaoka, Suita, Osaka Prefecture 565-0871, Japan

Table 1 Number of cycles for decoding for [6,14] and this paper

References	<i>m</i>	<i>n</i>	<i>t</i>	Bytes	sec	perm	synd	key eq	root	All	arch
[6]	13	6624	115	958482	252	23140	83127	102337	65050	444971	IB
	13	6960	119	1046739	263	23020	83735	109805	66453	456292	IB
[14]	13	8192	128	1357824	297	3783	62170	170576	53825	410132	IB
						3444	36076	127070	34491	275092	HW
This	13	8192	128	1357824	297	4853	62170	170576	53825	412272	IB
						4457	36076	127070	34491	277118	HW

The numbers in bold are copied from the rows for [14] as the implementations are the same

secret key: in order to generate uniformly random secret keys among the key space, we use a more sophisticated key generation routine for this paper. This issue is discussed in detail in Sect. 6. Also, in the decoding routine, the “permutation” is implemented in a slightly different way from that for [14]; see discussions below and Sect. 3.

Also note that the implementation presented in this paper is adapted from the SSE and AVX implementations for the parameter set *mceliece8192128* of the “Classic McEliece” submission [5] (as an IND-CCA2-secure KEM) to NIST’s call. The implementation for this paper is slightly more optimized in key-pair generation and decapsulation, and therefore, the cycle counts in this paper are better than what is shown in the submission.

1.1 Performance

The decoding timings of our software, the implementation for [14], as well as those for the highest-security parameters in [6, Table 1], are listed in Table 1. Most notations here are the same as in [6, Table 1]: we use *m* to indicate the field size 2^m , *n* to denote the code length, and *t* to denote the number of errors. “Bytes” is the size of public keys in bytes; “Sec” is the (pre-quantum) security level reported by the <https://bitbucket.org/cbcrypto/isdfq> script from Peters [24], rounded to the nearest integer. We list the cycle counts for each stage of the decoding process as in [6, Table 1]: “perm” for secret permutation, “synd” for syndrome computation, “key eq” for key-equation solving, and “root” for root finding. In [6, Table 1], there are two columns for “perm”: one stands for the initial permutation and one stands for the final permutation, but the cycle counts are essentially the same (we pick the timing for the initial permutation). Note that the column “all,” which serves as an estimation for the KEM decryption time, is computed as

$$\text{“perm”} \times 2 + \text{“synd”} \times 2 + \text{“key eq”} + \text{“root”} \times 2.$$

This is different from the “total” column in [6, Table 1] for decoding time, which is essentially

$$\text{“perm”} \times 2 + \text{“synd”} + \text{“key eq”} + \text{“root”}.$$

Table 2 Cycle counts for key-pair generation (key gen.), encapsulation (encap.), and decapsulation (decap.)

References	key gen.	encap.	decap.	arch
[14]	1552717680	312135	492404	IB
	1236054840	289152	343344	HW
This paper	3193838257	330839	523690	IB
	1578541256	282372	355152	HW
BIKE [1]	≈ 690000	≈ 360000	≈ 8270000	KL
KYBER [2]	121056	157964	154952	HW

The difference is explained in Sect. 6 in detail. “Arch” indicates the microarchitecture of the platform: “IB” for Ivy Bridge, “HW” for Haswell, and “KL” (in Table 2) for Kaby Lake.

We comment that the way we exploit internal parallelism brings some overhead that can be avoided when using external parallelism. In general, such an overhead is hard to avoid since the data flow of the algorithm is not necessarily friendly for bitslicing internally. This is exactly the main reason why our software is slower in “key eq” than McBits (a minor reason is that we are using a larger *t*). Despite the extra overhead, we still perform better when it comes to “synd” and “root.” The improvement on “perm” is mainly because of our use of an asymptotically faster algorithm. Our “all” speed ends up being better than McBits. We emphasize that the timings for McBits are actually 1/256 of the timings for 256 parallel decryption operations, while the timings for our software involve only one decryption operation.

For completeness, we also show the cycle counts of our implementation for key generation, encapsulation, and decapsulation in Table 2 and compare with those in [14]. All the cycle counts for our software were measured using the SUPERCOP benchmarking toolkit. We use the `randombytes` function in SUPERCOP to obtain randomness for the key generation and encapsulation. In recent versions of SUPERCOP, `randombytes` first reads from the operating system entropy to obtain a short seed, and then, the seed is expanded using some pseudorandom number generator to generate a byte string of the desired length. According

Table 3 Numbers of bytes for secrets keys (sk. size), public keys (pk. size), and the ciphertexts (ct. size) for the submissions to NIST’s call

References	sk. size	pk. size	ct. size
This paper [5]	14080	1357824	240
BIKE [1]	548	8188	8188
KYBER [2]	3168	1440	1504

to our experiment results, the calls to `randombytes` are not significant in either key generation or encapsulation.

For comparison with other submissions to NIST’s call, Table 2 also includes the speeds for the IND-CPA-secure code-based KEM “BIKE” [1] and the IND-CCA2-secure lattice-based KEM “KYBER” [2]. We use the largest parameter set of the variant BIKE-1 for BIKE and the parameter set KYBER1024 for KYBER. All the parameter sets (including the one for our implementation) are designed to fit the level-5 security level defined by NIST. All implementations being compared are constant time. The sizes of the secret keys, public keys, and the ciphertexts for the submission are summarized in Table 3.

1.2 Parameter selection

As shown in Table 1, we implement one specific parameter set

$$(m, n, t) = (13, 8192, 128),$$

with 1357824-byte public keys and a 2^{297} security level. We explain below the reasons to select this parameter set.

The Gao–Mateer additive FFT evaluates the input polynomial at a predefined \mathbb{F}_2 -linear subspace of \mathbb{F}_{2^m} . The parameter n indicates the size of the list of field elements that we need to evaluate at, so for $n = 2^m$ we can simply define the subspace as \mathbb{F}_{2^m} . In the case of $n < 2^m$, however, there is no way to define the subspace to fit arbitrary choice of the field elements (which is actually a part of the secret key), so the best we can do is still evaluate at the whole \mathbb{F}_{2^m} . In other words, having $n < 2^m$ would result in some redundant computation.

The parameter n also indicates the number of elements that we need to apply secret permutations on. The permutation algorithm we use, in its original form, requires that the number of elements to be a power of 2. The algorithm can be “truncated” to deal with an arbitrary number of elements, but this makes implementation difficult.

Having t close to the register size is convenient for bit-slicing the FFT algorithms and also the Berlekamp–Massey algorithm. We choose $t = 128$ to match the size of XMM registers in SSE-supporting architectures, as well as the size of the vector registers in the ARM-NEON architectures. Not having t close to the register size will not really affect the

performance of FFTs: the algorithms are dominated by the t -irrelevant part as long as t is much smaller than 2^m . A bad value for t has more impact on the performance of the Berlekamp–Massey algorithm since we might waste many bits in the registers. Choosing $t = 128$ (after choosing $n = 2^m$) also forces the number of rows mt and number of columns $n - mt$ of the public-key matrix to be multiples of 128, which is convenient for implementing the encryption operation.

For the reasons stated above, some other nice parameters for (m, n, t) are

- (12, 4096, 64) with 319488-byte public keys and a 2^{159} security level,
- (12, 4096, 128) with 491520-byte public keys and a 2^{189} security level, and
- (13, 8192, 64) with 765440-byte public keys and a 2^{210} security level.

We decided to select a parameter set that achieves at least a 2^{256} pre-quantum security level and thus presumably at least a 2^{128} post-quantum security level.

The reader might argue that such a high security level is not required for real applications. Indeed, even if quantum algorithms can take a square root on the security level, it still means that our system has a roughly 2^{150} post-quantum security level. In fact, we even believe that quantum algorithms will not be able to take a square root on the security: we believe there is an overhead of more than 2^{20} that needs to be added upon the square root. However, before the post-quantum security of our system is carefully analyzed, we think it is not a bad idea to implement a parameter set that is very likely to be an overkill and convince users that the system achieves a decent speed even in this case. Once careful analysis is done, our implementation can then be truncated to fit the parameters. The resulting implementation will have at least the same speed and a smaller key size.

1.3 Software availability

Our software is available at <https://tungchou.github.io/code/mceliece8192128.tar.gz>. The code can be easily put into SUPERCOP for benchmarking; see <https://bench.cr.yp.to/tips.html>. The whole software is in the public domain.

1.4 Organization

The rest of this paper is organized as follows. Section 2 introduces the low-level building blocks used in our software. Section 3 describes how we implement the Beneš networks for secret permutations. Section 4 describes how we implement the Gao–Mateer FFT for root finding and the

corresponding “transposed” FFT for syndrome computation. Section 5 introduces how we implement the Berlekamp–Massey algorithm for key-equation solving. Finally, Sect. 6 introduces how the components in Sects. 3, 4 and 5 are combined to form the complete decryption, as well as how key generation and encryption are implemented.

2 Building blocks

This section describes the low-level building blocks used in our software. We will use these building blocks as black boxes in the following sections. The implementation techniques behind these building blocks are not new. In particular, this section presents (1) how to use bitslicing to perform several field operations in parallel and (2) how to perform bit-matrix transposition in software. Readers who are familiar with these techniques may skip this section.

2.1 Individual field operations

The finite field $\mathbb{F}_{2^{13}}$ is constructed as $\mathbb{F}_2[x]/(g)$, where $g = x^{13} + x^4 + x^3 + x + 1$. Let $z = x + (g)$. Each field element $\sum_{i=0}^{12} a_i z^i$ can then be represented as the integer $(a_{12}a_{11} \cdots a_0)_2$ in software. Field additions are carried out by XORs between integers. Field multiplications are carried out by the following C function.

```
typedef uint16_t gf;
gf gf_mul(gf in0, gf in1)
{
    uint64_t i, tmp, t0=in0, t1=in1, t;
    tmp = t0 * (t1 & 1);
    for (i = 1; i < 13; i++)
        tmp ^= (t0 * (t1 & (1 << i)));
    t = tmp & 0x1FF0000;
    tmp ^= (t >> 9)^(t >> 10)^(t >> 12)^(t >> 13);
    t = tmp & 0x000E000;
    tmp ^= (t >> 9)^(t >> 10)^(t >> 12)^(t >> 13);
    return tmp & ((1 << 13)-1);
}
```

The squaring function is written in a similar way. Computing the inverse of a field element is carried out by raising the element to the power $2^{13}-2$ using 12 squarings and 4 multiplications.

2.2 Bitsliced field operations

The field multiplication function `gf_mul` and the field addition shown above are rather inefficient. The reason is that each logical instruction deals with only a small number of bits. For the algorithms used in our software, however, most of the time several field operations can be performed in parallel. We thus “bitslice” the field operations. The idea of bitslicing is to use bitwise logical operations to simulate w copies of a combinational circuit: the data for the i th copy are stored

```
void vec64_mul(uint64_t *h, uint64_t *f, uint64_t *g)
{
    int i, j;
    uint64_t r[2*13 - 1];
    for (i = 0; i < 2*13 - 1; i++)
        r[i] = 0;
    for (i = 0; i < 13; i++)
        for (j = 0; j < 13; j++)
            r[i+j] ^= (f[i] & g[j]);
    for (i = 2*13-2; i >= 13; i--)
    {
        r[i - 9] ^= r[i];
        r[i - 10] ^= r[i];
        r[i - 12] ^= r[i];
        r[i - 13] ^= r[i];
    }
    for (i = 0; i < 13; i++)
        h[i] = r[i];
}
```

Fig. 1 The C function for bitsliced multiplications in $\mathbb{F}_{2^{13}}[x]/(x^{13} + x^4 + x^3 + x + 1)$ using 64-bit words

in the i th bits of the registers. In this way, the number of bits involved in each instruction can be improved to w . Bitslicing is also heavily used in [6]. We emphasize that for [6], the w copies are from w different decryption operations. For our software, the w copies are all from the same decryption operation.

The function `vec64_mul` for bitsliced field multiplications using 64-bit words is shown in Fig. 1. Our software uses 128-bit or 256-bit words instead. According to Fog’s well-known performance survey [16], on the Ivy Bridge architecture, the bitwise AND/XOR/OR instructions on the 128-bit registers (XMM registers) have a throughput of 3 per cycle, while for the 256-bit registers (YMM registers) the throughput is only 1. On Haswell, the instructions for the 256-bit registers have a throughput of 3 per cycle. We thus use the corresponding function `vec128_mul` for Ivy Bridge and use `vec256_mul` as much as possible for Haswell. Since both functions are heavily used in our software, they are written in `qhasm` [4] code for the best performance.

Many CPUs nowadays support the `pclmulqdq` instruction. The instruction essentially performs a multiplication between two 64-coefficient polynomials in $\mathbb{F}_2[x]$, so it can be used for field multiplications. Our multiplication function `vec256_mul` takes 138 Haswell cycles, which means a throughput of 1.86 field multiplications per cycle. The `pclmulqdq` instruction has a throughput of 1/2 on Haswell. We may perform 2 multiplications between 13-coefficient polynomials using one `pclmulqdq` instruction. However, non-bitsliced representations make it expensive to perform reductions modulo the irreducible polynomial g . The throughput for `pclmulqdq` is only 1/8 on Ivy Bridge, which makes it even less favorable.

2.3 Transposing bit matrices

Bit-matrix transposition appears to be a well-known technique in computer programming. Perhaps due to the simplicity of the method, it is hard to trace who the credit belongs to. Below, we give a brief review on the idea.

The task is to transpose a $w \times w$ bit matrix M , where w is a power of 2. The idea is to first divide the matrix into $4 \times w/2 \times w/2$ submatrices, i.e., the left upper, right upper, left bottom, and right bottom submatrices. Then, a “coarse-grained transposition” is performed on M , which simply interchanges the left bottom and right upper submatrices. Finally, each block is transposed recursively, until we reach 1×1 matrices. The idea is depicted below.

$$M = \begin{pmatrix} M_{00} & M_{01} \\ M_{10} & M_{11} \end{pmatrix} \implies \begin{pmatrix} M_{00} & M_{10} \\ M_{01} & M_{11} \end{pmatrix} = M'$$

$$\implies \begin{pmatrix} M_{00}^T & M_{10}^T \\ M_{01}^T & M_{11}^T \end{pmatrix} = M^T$$

The benefit of this approach is that it can be carried out efficiently in software. Suppose we are working on a w -bit machine, where the matrix is naturally represented as an array of w -bit words in a row-major fashion. Observe that each of the first $w/2$ rows of M' is the concatenation of the first halves of two rows in M . Similarly, each of the second $w/2$ rows is the concatenation of the second halves of two rows in M . Therefore, each row in M' can be generated using a few logical operations. After this, in order to carry out operations in the recursive calls efficiently, the operations involving the upper two blocks can be handled together using logical operations on w -bit words. The same applies for the bottom two blocks. The C code for transposing 64×64 matrices is shown in Fig. 2.

The same technique can be easily generalized to deal with non-square matrices. Our software makes use of functions for transposing 64×128 and 128×64 matrices, where instructions such as `psrlq`, `psllq`, `psrld`, `pslld`, `psrlw`, and `psllw` are used to shift the 128-bit registers.

3 The Beneš network

As described in [6], a “permutation network” uses a sequence of conditional swaps to apply an arbitrary permutation to an input array S . Each conditional swap is a permutation-independent pair of indices (i, j) together with a permutation-dependent bit c ; it swaps $S[i]$ with $S[j]$ if $c = 1$. Our software uses a specific type of permutation network, called the Beneš network [8], to perform secret permutations for the code-based encryption system.

```
const uint64_t m[6][2] =
{
  {0x5555555555555555, 0xAAAAAAAAAAAAAAAA},
  {0x3333333333333333, 0xcccccccccccccccc},
  {0xf0f0f0f0f0f0f0f0, 0xf0f0f0f0f0f0f0f0},
  {0x00ff00ff00ff00ff, 0xff00ff00ff00ff00},
  {0x0000ffff0000ffff, 0xffff0000ffff0000},
  {0x00000000ffff0000, 0xffff000000000000}
};
for (j = 5; j >= 0; j--)
{
  s = 1 << j;
  for (p = 0; p < 32/s; p++)
  for (i = 0; i < s; i++)
  {
    idx0 = p*2*s + i;
    idx1 = p*2*s + i + s;
    x = (in[idx0] & m[j][0]) | ((in[idx1] & m[j][0]) << s);
    y = ((in[idx0] & m[j][1]) >> s) | (in[idx1] & m[j][1]);
    in[idx0] = x;
    in[idx1] = y;
  }
}
```

Fig. 2 The C code for transposing 64×64 bit matrices. The matrix to be transposed is stored in the array `in`. The transposition is performed in-place

The McBits paper uses a “sorting network” for the same purpose but notes that it takes asymptotically more conditional swaps than the Beneš network: $O(n \log^2 n)$ versus $O(n \log n)$ for array size $n = 2^m$. We found that the Beneš network is more favorable for our implementation because it is easier to use the internal parallelism due to its simple structure. This section introduces the structure of the Beneš network, as well as how it is implemented in our software.

3.1 Conditional swaps: structure

The Beneš network for 2^m elements consists of a sequence of $2m - 1$ stages, where each stage consists of exactly 2^{m-1} conditional swaps. The set of index pairs for these 2^{m-1} conditional swaps is defined as

$$\{(\alpha \cdot 2^{s+1} + \beta, \alpha \cdot 2^{s+1} + 2^s + \beta) \mid 0 \leq \alpha < 2^{m-1-s}, 0 \leq \beta < 2^s\},$$

where s is stage dependent. The sequence of s is defined as

$$0, 1, \dots, m - 2, m - 1, m - 2, \dots, 1, 0.$$

To visualize the structure, the size-16 Beneš network is depicted in Fig. 3.

The Beneš network is often defined in a recursive way, in which case the size- 2^m Beneš network is viewed as the combination of the first and last stage, plus 2 size- 2^{m-1} Beneš networks in the middle.

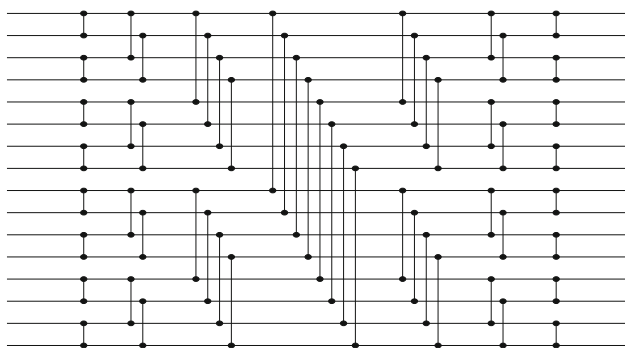


Fig. 3 The size-16 Beneš network with 7 stages. Each horizontal line represents an element in the array. Each vertical line segment illustrates a conditional swap involving the array elements corresponding to the end points

3.2 Conditional swaps: implementation

Consider the Beneš network for an array S of 2^m bits where $m = 13$. We consider S as a 128×64 matrix M such that

$$M_{i,j} = S[i + 128j].$$

In each of the first and last 7 stages, the index pairs always have an index difference that is in

$$\{1, 2, 4, 8, 16, 32, 64\}.$$

This implies that in each of these stages, $M_{i,j}$ is always conditionally swapped with $M_{i',j}$, where i' depends only on i . This implies that the conditional swaps can be carried out by performing bitwise logical operations between the rows (and the vectors formed by the corresponding conditions): a conditional swap between $M_{i,j}$ and $M_{i',j}$ with condition bit c can be carried out by 4-bit operations

$$\begin{aligned} y &\leftarrow M_{i,j} \oplus M_{i',j}; \\ y &\leftarrow cy; \\ M_{i,j} &\leftarrow M_{i,j} \oplus y; \\ M_{i',j} &\leftarrow M_{i',j} \oplus y; \end{aligned}$$

as mentioned in [6]. Likewise, the 11 stages in the middle can be carried out by using bitwise logical operations between columns.

Because of the relationship between M and S , in our implementation M is naturally stored in a column-major fashion at the beginning. In order to perform the operations between the rows in the first 7 stages, we transpose the matrix (using the technique described in Sect. 2) to obtain the row-major representation. In order to perform the operations between the columns in the middle stages, another transposition is performed to obtain the column-major representation. Similarly, another transposition is performed at the end of the

middle stages, and the last transposition is performed after all the stages.

Note that there are 4 transpositions in total. In the implementation for [14], the first and last transpositions are omitted. Omitting the transpositions are fine, as all we need is to redefine M and S . For this paper, we decide not to omit the transposition to make the relationship easy to define.

For the middle stages, we store each column in a 128-bit word, and the operations are carried out using bitwise logical instructions on XMM registers. For the first and last 7 stages, each row is merely 64 bits, so a straightforward way is to perform bitwise logical instructions on the general-purpose 64-bit registers. We do better by combining two 64-bit logical operations into one logical operation on XMM registers.

4 The Gao–Mateer additive FFT

Given a predefined \mathbb{F}_2 -linear basis $\{\beta_1, \beta_2, \dots, \beta_k\} \subset \mathbb{F}_{2^m}$ and an ℓ -coefficient input polynomial

$$f = \sum_{i=0}^{\ell-1} f_i x^i \in \mathbb{F}_{2^m}[x]$$

such that $\ell \leq 2^k \leq 2^m$, the Gao–Mateer FFT evaluates f at all the subset sums of the basis. In other words, the FFT outputs the sequence

$$f(e_1), f(e_2), \dots, f(e_{2^k}),$$

where

$$(e_1, e_2, e_3, e_4, e_5, \dots) = (0, \beta_1, \beta_2, \beta_1 + \beta_2, \beta_3, \dots).$$

Such an FFT will be called a size- 2^k FFT.

Assuming that $\beta_k = 1$. The idea is to compute two polynomials $f^{(0)}$ and $f^{(1)}$ such that

$$f = f^{(0)}(x^2 + x) + x f^{(1)}(x^2 + x),$$

using the “radix conversion” described in [6, Section 3] (this is called “Taylor expansion” in [17]). Note that $f^{(0)}$ is a $\lceil \ell/2 \rceil$ -coefficient polynomial, while $f^{(1)}$ is a $\lfloor \ell/2 \rfloor$ -coefficient polynomial. Observe that $\alpha^2 + \alpha = (\alpha + 1)^2 + (\alpha + 1)$. This implies that once $t_0 = f^{(0)}(\alpha^2 + \alpha)$ and $t_1 = f^{(1)}(\alpha^2 + \alpha)$ are computed, $f(\alpha)$ can be computed as $t_0 + \alpha \cdot t_1$, and $f(\alpha + 1)$ can be computed as $f(\alpha) + t_1$. Observe that the output of the FFT is the sequence

$$\begin{aligned} &f(e_1), f(e_2), \dots, f(e_{2^{k-1}}), \\ &f(e_1 + 1), f(e_2 + 1), \dots, f(e_{2^{k-1}} + 1), \end{aligned}$$

and $e_1, \dots, e_{2^{k-1}}$ form all subset sums of $\{\beta_1, \dots, \beta_{k-1}\}$.

Therefore, two FFT recursive calls are carried out to evaluate $f^{(0)}$ and $f^{(1)}$ at all subset sums of

$$\{\beta_1^2 + \beta_1, \dots, \beta_{k-1}^2 + \beta_{k-1}\}.$$

Finally, $f(e_i)$ and $f(e_i + 1)$ are computed by using $f^{(0)}(e_i^2 + e_i)$ and $f^{(1)}(e_i^2 + e_i)$ from the recursive calls, for all i from 1 to 2^{k-1} .

In the case where $\beta_k \neq 1$, the task is reconsidered as evaluating $f(\beta_k x)$ at the subset sums of

$$\{\beta_1/\beta_k, \beta_2/\beta_k, \dots, 1\}.$$

This is called “twisting” in [6]. Note that it takes $\ell - 1$ multiplications to compute $f(\beta_k x)$. To summarize, the Gao–Mateer additive FFT consists of 4 steps: (1) twisting, (2) radix conversion, (3) two FFT recursive calls, and (4) combining outputs from the recursive calls.

In order to find the roots of an error locator, we need to evaluate at every field element in $\mathbb{F}_{2^{13}}$. The corresponding basis is defined as

$$\{\beta_1 = z^{12}, \beta_2 = z^{11}, \dots, \beta_{13} = 1\}.$$

Having $\beta_{13} = 1$ means that the first twisting can be skipped. Since we use $t = 128$, the error locator for our system is a 129-coefficient polynomial. However, for implementation of the FFT algorithm it is more convenient to have a 128-coefficient input polynomial. We therefore consider the error locator as $x^{128} + f$ and compute $\alpha^{128} + f(\alpha)$ for all $\alpha \in \mathbb{F}_{2^{13}}$. Below, we explain how the Gao–Mateer additive FFT for root finding, as well as the corresponding “transposed” FFT for syndrome computation, are implemented in our software.

4.1 Radix conversions and twisting

In [6], it is described that the first step of the radix conversion is to compute polynomials Q and R from the $4n$ -coefficient (n is a power of 2) input polynomial $f = \sum_{i=0}^{4n-1} f_i x^i$:

$$Q = (f_{2n} + f_{3n}) + \dots + (f_{3n-1} + f_{4n-1})x^{n-1} + f_{3n}x^n + \dots + f_{4n-1}x^{2n-1},$$

$$R = (f_0) + \dots + (f_{n-1})x^{n-1} + (f_n + f_{2n} + f_{3n})x^n + \dots + (f_{2n-1} + f_{3n-1} + f_{4n-1})x^{2n-1},$$

so that $f = Q(x^{2n} + x^n) + R$. Then, Q and R are fed into recursive calls to obtain the corresponding $R^{(0)}, R^{(1)}, Q^{(0)}, Q^{(1)}$. Finally, the routine outputs $f^{(0)} = R^{(0)} + x^n Q^{(0)}$

and $f^{(1)} = R^{(1)} + x^n Q^{(1)}$. The recursion ends when we reach a 2-coefficient polynomial $f_0 + f_1 x$, in which case $f^{(0)} = f_0$ and $f^{(1)} = f_1$.

Here is a straightforward way to implement the routine. First of all, represent the input polynomial f as a $4n$ -element array `in` of datatype `gf` (see Sect. 2) such that f_i is stored in `in[i]`. Then, perform $4n$ XORs

```
for (i = 0; i < n; i++) in[2*n+i] ^= in[3*n+i];
for (i = 0; i < n; i++) in[1*n+i] ^= in[2*n+i];
```

to store R_i in `in[i]` and Q_i in `in[2*n+i]`. Likewise, the additions in the recursive calls can be carried out by in-place XORs between array elements. Eventually, we have $f_i^{(0)}$ in `in[2*i]` and $f_i^{(1)}$ in `in[2*i+1]`.

Representing the polynomials as arrays in `gf` is, however, expensive for twisting: as mentioned in Sect. 2, the function `gf_mul` is not efficient. Therefore in our software, the polynomials are represented in bitsliced format. In this case, the additions can be simulated by using bitwise logical instructions and shifts. As a concrete example, let f be a 64-coefficient input polynomial in $\mathbb{F}_{2^{13}}[x]$, which is represented as a 13-element array of type `uint64_t`. Then, the following code applies the radix conversion on f .

```
const uint64_t mask[5][2] =
{
    {0x8888888888888888, 0x4444444444444444},
    {0xC0C0C0C0C0C0C0C0, 0x3030303030303030},
    {0xF000F000F000F000, 0x0F000F000F000F00},
    {0xFF000000FF000000, 0x00FF000000FF0000},
    {0xFFFF000000000000, 0x0000FFFF00000000}
};
for (k = 4; k >= 0; k--)
for (i = 0; i < 13; i++)
{
    in[i] ^= (in[i] & mask[k][0]) >> (1 << k);
    in[i] ^= (in[i] & mask[k][1]) >> (1 << k);
}
```

In the end, the coefficients of $f^{(0)}$ are represented by the even bits of the words, while the coefficients of $f^{(1)}$ are represented by the odd bits.

The same technique can also be used to complete the radix conversions in the FFT recursive calls. Since a twisting operation simply multiplies f_i by β_k^i , they are carried out using bitsliced multiplications. See Fig. 4 for the code for all the radix conversions and twisting operations, including those in the FFT recursive calls. Note that the first twisting operation, which should take place before the first radix conversion, is already skipped in the code. Our software uses similar code but replaces 64-bit words by 128-bit words.

4.2 Butterflies

The reader might have noticed that the last 4 stages of Fig. 3 are similar to the well-known butterfly diagram for standard

multiplicative FFTs. In a standard multiplicative FFT, f is written as $f^{(0)}(x^2) + xf^{(1)}(x^2)$ so that $f(\alpha)$ and $f(-\alpha)$ can be computed using $f^{(0)}(\alpha^2)$ and $f^{(1)}(\alpha^2)$ obtained from recursive calls. The similarity (between multiplicative FFTs and additive FFTs) in the ways of rewriting f results in the same “butterfly” structure.

In the case of a “full-size” additive FFT, where $\ell = 2^k$, the whole butterfly diagram has to be carried out. The technique used for carrying out the Beneš network (see Sect. 3) can be easily generalized to carry out the diagram. For decoding, however, ℓ is usually much smaller than $2^k = 2^m$. As the result, we only need to carry out the last $\log_2 \ell$ stages of the complete butterfly diagram.

As described in Sect. 3, we carry out the second half of the Beneš network by using a bit-matrix transposition in the middle. In the case of additive FFT butterflies, there will be m bit-matrix transpositions. The ideal case is that the ℓ is small enough so that the transpositions can be avoided. The corresponding code using 64-bit words for $m = 12$ is presented in Fig. 5. For the parameters $\ell = 128$ and $m = 13$, we are close to this ideal case but need to carry out 1 or 2 extra stages. The extra stages can be carried out by interleaving the 128-bit or 256-bit words.

4.3 The bottom level of recursion

As shown in Fig. 4, when carrying out the radix conversions and twisting operations, we maintain a list of ℓ field elements. On the other hand, as shown in Fig. 5, when carrying out the FFT butterflies, we maintain a list of 2^m field elements. Apparently, some operations are required to transit from the ℓ -element representation to the 2^m -element representation. This has to do with how the bottom level of recursion is defined.

The straightforward way to end the recursion is to check whether the input polynomial has only 1 coefficient; if so, the output is simply copies of the coefficient (the constant term). This is exactly the case for Figs. 4 and 5: after running the code in Fig. 4, we simply prepare the bitsliced representation of 64 copies of each element and store them in `out`, and then, Fig. 5 can be run to complete the FFT.

We do better by using the idea in [6, Section 3] to end the recursion when the input is a 2-coefficient polynomial. Let the input be $f = f_0 + f_1x$ and the basis be $\{\beta_1, \dots, \beta_k\}$. The idea is to first prepare a table containing $f_1\beta_i$ for all i , and then, each output element can be computed using at most one field addition. To implement the idea, we perform the radix conversions and twisting operations as in Fig. 4 but stop when we reach 2-coefficient polynomials. At this moment, the $\ell/2$ elements corresponding to f_0 would lie in the lower $\ell/2$ bits of the ℓ -bit words, while those for f_1 would lie in the higher $\ell/2$ bits. The outputs of the lowest-level FFTs can then be obtained by carrying out bitsliced multiplications

```
for (j = 0; j <= 4; j++)
{
  for (i = 0; i < 13; i++)
  for (k = 4; k >= j; k--)
  {
    in[i] ^= (in[i] & mask[k][0]) >> (1 << k);
    in[i] ^= (in[i] & mask[k][1]) >> (1 << k);
  }
  vec64_mul(in, in, s[j]); // twisting
}
```

Fig. 4 The code for performing the twisting operations and radix conversion in the FFT for a 64-coefficient polynomial $f \in \mathbb{F}_{2^{13}}[x]$

```
for (i = 0; i <= 5; i++)
{
  s = 1 << i;
  for (j = 0; j < 64; j += 2*s)
  for (k = j; k < j+s; k++)
  {
    vec64_mul(tmp, out[k+s], consts[ptr + (k-j)]);
    for (b = 0; b < 13; b++)
      out[k][b] ^= tmp[b];
    for (b = 0; b < 13; b++)
      out[k+s][b] ^= out[k][b];
  }
  ptr += (1 << i);
}
```

Fig. 5 Butterflies in the additive FFT

and additions using bitwise logical operations between the $\ell/2$ -bit words.

After this, we have the bitsliced representation (an array of m $\ell/2$ -bit words) for the first output elements of the lowest-level FFTs, the representation for the second output elements, and so on; in total there are $2^m/(\ell/2)$ such arrays. In order to group the output elements that belong to the same lowest-level FFT, we perform a sequence of m transpositions on $2^m/(\ell/2) \times (\ell/2) = 128 \times 64$ bit matrices, using the technique described in Sect. 2. Finally, the FFT butterflies can be performed using code similar to Fig. 5.

4.4 The transposed additive FFT

As described in [6, Section 4], a *linear algorithm* can be represented as a directed graph, and an algorithm that performs the transposed linear map can be obtained by reversing the edges in the graph. The way we implement the FFT makes it easy to imagine the structure of the graph and program the corresponding transposed FFT. As shown in Figs. 4 and 5, each inner loop in our FFT code essentially applies a simple linear operation on the values in `in` or `out`. In general, it suffices to modify the loops to reverse the order that the inner loop is iterated and then replace


```

for (i = 5; i >= 0; i--)
{
    s = 1 << i;
    ptr -= s;
    for (j = 0; j < 64; j += 2*s)
    for (k = j; k < j+s; k++)
    {
        for (b = 0; b < 13; b++) out[k][b] ^= out[k+s][b];
        vec64_mul(tmp, out[k], consts[ ptr + (k-j) ]);
        for (b = 0; b < 13; b++) out[k+s][b] ^= tmp[b];
    }
}

:
:
:
for (j = 4; j >= 0; j--)
{
    vec64_mul(in, in, s[j]); // twisting
    for (k = j; k <= 4; k++)
    for (i = 0; i < 13; i++)
    {
        in[i] ^= (in[i] & (mask[k][1] >> (1 << k)))
                << (1 << k);
        in[i] ^= (in[i] & (mask[k][0] >> (1 << k)))
                << (1 << k);
    }
}
    
```

Fig. 6 Transposed FFT code with respect to Figs. 4 and 5

the inner loop by its transpose. The transposed additive FFT code with respect to Figs. 4 and 5 is shown in Fig. 6 (the code for transposing the bottom level of recursion is skipped).

5 The Berlekamp–Massey algorithm

The description of the original Berlekamp–Massey algorithm (BM) can be found in [19]. In each iteration of the algorithm, a field inversion has to be carried out. To perform the inversion in constant time, we may use the square-and-multiply algorithm, but this is rather expensive as discussed in Sect. 2. To avoid the problem, our implementation follows the inversion-free version of the algorithm as described in [27].

The algorithm begins with initializing polynomials $\sigma(x) = 1, \beta(x) = x \in \mathbb{F}_{2^m}[x], \ell = 0 \in \mathbb{Z},$ and $\delta = 1 \in \mathbb{F}_{2^m}.$ The input syndrome polynomial is denoted as $S(x) = \sum_{i=0}^{2t-1} S_i x^i.$ Then in iteration k (from 0 to $2t - 1$), the variables are updated using operations in Fig. 7. Note that ℓ and δ are just an integer and a field element, and multiplying a polynomial by x (to update $\beta(x)$) is rather cheap. Therefore, the algorithm is bottlenecked by computing d and updating $\sigma(x).$ We explain below how the algorithm is implemented in our software.

$$d \leftarrow \sum_{i=0}^t \sigma_i S_{k-i}$$

$$[\sigma(x), \beta(x), \ell, \delta] \leftarrow \begin{cases} [\delta\sigma(x) - d\beta(x), x\beta(x), \ell, \delta], \\ d = 0 \text{ or } k < 2\ell. \\ [\delta\sigma(x) - d\beta(x), x\sigma(x), k - \ell + 1, d], \\ \text{otherwise.} \end{cases}$$

Fig. 7 Iteration k in the inversion-free BM

5.1 General implementation strategy

Assume that there are $(t + 1)$ -bit general-purpose registers on the target machine. For example, one can assume that $t = 63$ and that we are working on a 64-bit machine. We store polynomials $\sigma(x)$ and $\beta(x)$ in the bitsliced format, each using an array of $m(t + 1)$ -bit words. The constant terms σ_0 and β_0 are stored in the most significant bits of the words; σ_1 and β_1 are stored in the second significant bits; and so on. We also use an array S' of $m(t + 1)$ -bit words to store at most $t + 1$ coefficients of $S(x).$ This array is maintained so that S_k is stored in the most significant bits of the words; S_{k-1} is stored in the second significant bits; and so on.

To compute $d,$ we first perform a bitsliced field multiplication between $\sigma(x)$ and $S'.$ The result is the bitsliced representation of $\sigma_0 S_k, \sigma_1 S_{k-1},$ etc. The element d can then be computed as the parities of the $m(t + 1)$ -bit words. After this, S_{k+1} is inserted to the most significant bits of the words in $S',$ which will be used in the next iteration.

To update $\sigma(x),$ we need to perform two scalar multiplications $\delta \cdot \sigma(x)$ and $d \cdot \beta(x).$ The bitsliced representations of $t + 1$ copies of δ and d are first prepared, and then, bitsliced multiplications are carried out to compute the products. Updating $\beta(x)$ is done by conditionally replacing the value of $\beta(x)$ by $\sigma(x)$ (which can be easily represented as logical operations) and then shifting each word to the right by one bit to simulate the multiplication by $x.$

The implementation strategy pretty much simulates the circuit presented in [27, Figure 1]. Using the strategy, (each iteration of) the BM algorithm can be represented as a fixed sequence of instructions. In particular, the load and store instructions always use the same memory indices. As a result, the implementation is fully protected against timing attacks.

5.2 Haswell Implementation for $t = 128$

Exactly, the same implementation strategy cannot be used for $t = 128$ on Haswell for there are no $(128 + 1)$ -bit registers. To solve this problem, our strategy is to store σ_0 and S_k in

two variables of datatype gf . The elements $\sigma_1, \dots, \sigma_{128}$ and S_{k-1}, \dots, S_0 are still stored in the bitsliced format, using two arrays of 128-bit words. To compute d , the product $\sigma_0 S_k$ is computed separately. Similarly, to update $\sigma(x)$, the product $\sigma_0 \delta$ is computed separately. Note that β_0 is always 0, so we simply store $\beta_1, \dots, \beta_{128}$ in the bitsliced format.

We also need a way to update S' and $\beta(x)$ without generic shift instructions for 128-bit registers. Our solution is to make use of the `shrd` instruction. Given 64-bit registers r_1, r_0 as arguments, the `shrd` instruction is able to shift the least significant bit of r_1 into the most significant bit of r_0 . Therefore, with 2 `shrd` instructions, we can shift a 128-bit word by one bit to the right. In particular, the second `shrd` shifts one bit into the most significant bit of the 128-bit word. Therefore, we update S' by setting this bit to bits of S_k and update β by setting this bit to 0 or bits of σ_0 (depending on the condition).

To optimize the speed for Haswell, we combine the two `vec128_mul` function calls for $\delta \cdot \sigma(x)$ and $d \cdot \beta(x)$ to form one `vec256_mul`. As discussed in Sect. 2, this is better because 256-bit logical instructions have the same throughput as the 128-bit ones.

We also use 256-bit logical instructions to accelerate `vec128_mul`. A field multiplication can be viewed as a multiplication between 13-coefficient polynomials, followed by a reduction modulo g . Let the polynomials be f and f' ; the idea is to split the polynomial multiplication into two parts $f(f'_0 + \dots + f'_6 x^6)$ and $f(f'_7 + \dots + f'_{12} x^5 + 0x^6)$. In this way, we create two bitsliced multiplications for computing d , and the two can be combined as what we do for $\delta \cdot \sigma(x)$ and $d \cdot \beta(x)$. Note that for combining the two products and the reduction part we still use 128-bit logical instructions. By using 256-bit logical instructions, we improve the cycle counts of `vec128_mul` from 137 to 94 Haswell cycles.

As a minor optimization, we also combine the computation of $\sigma_0 S_k$ and $\sigma_0 \delta$. This is achieved by using the upper 32 bits of the 64-bit variables in `gf_mul` for another multiplication. In this way, two field multiplications can be carried out in roughly the same time as `gf_mul`.

As discussed in Sect. 4, the input of the FFT function for root finding is the bitsliced representation of f_0, \dots, f_{127} ; f_{128} is not stored since it is assumed to be 1. In fact, at the end of the Berlekamp–Massey algorithm we have $f_i = \sigma_{128-i}$. Therefore, we perform a field inversion for σ_0 and bitsliced multiplications to force a monic output polynomial for the Berlekamp–Massey algorithm.

6 The complete cryptosystem

The core of Classic McEliece, as in McBits, is essentially the Niederreiter cryptosystem. We briefly show below the how McBits and Classic McEliece achieve public-key encryption and key encapsulation using Niederreiter. For readers who are

interested in the detailed specification of Classic McEliece, please refer to [5]. We also show below how key-pair generation, encapsulation, and decapsulation are implemented in our software using the building blocks introduced in the previous sections, with focus on the Niederreiter-related operations.

6.1 McBits

In [6, Section 6], a complete public-key encryption system is described. The cryptosystem uses a KEM/DEM-like structure, where the KEM is based on the Niederreiter cryptosystem. To send a message, the sender first uses the receiver's Niederreiter public key to compute the syndrome of a random weight- t error vector. Then, the error vector is hashed to obtain two symmetric keys. The first symmetric key is used for a stream cipher to encrypt the arbitrary-length message. The second symmetric key is used for a message authentication code to authenticate the output generated by the stream cipher. The syndrome, the stream-cipher output, and the authentication tag are then sent to the receiver.

The receiver first decodes the syndrome using the Niederreiter secret key. The resulting error vector is then hashed to obtain the symmetric keys, and the receiver verifies (using the tag) and decrypts the stream-cipher output. Note that the receiver can fail in decoding or verification. The decryption algorithm should be carefully implemented such that others cannot distinguish (for example, by using timing information) what kind of failure the receiver encounters.

6.2 Classic McEliece

The specification of the keys and ciphertexts in Classic McEliece is a bit different from that of McBits:

- The secret key contains a random n -bit string s in addition to the Goppa polynomial and the support, and
- the ciphertext $C = (C_0, C_1)$ contains a hash value C_1 of the error vector e , in addition to the Niederreiter ciphertext C_0 .

C_1 serves as a “confirmation”: it is assumed that C_1 cannot be computed without knowledge of e .

During encapsulation, we compute $C_1 = \mathcal{H}(2, e)$ and $K = \mathcal{H}(1, e, C)$, where \mathcal{H} is a hash function. $C = (C_0, C_1)$ is then outputted as the ciphertext, and K is then outputted as the session key. During decapsulation, the received ciphertext C is first parsed as (C_0, C_1) . Then, a decoding algorithm is performed on C_0 ; the result of decoding is either an error vector e or \perp . In case of \perp , e is set to the random secret string s . We then compute $C'_1 = \mathcal{H}(2, e)$ and set e to s if $C'_1 \neq C_1$. A bit b is set to 0 if the decoding algorithm outputs

\perp or $C'_1 \neq C_1$; b is set to 1 otherwise. Finally, the session key K is set to be $\mathcal{H}(b, e, C)$.

6.3 Private-key generation

The private key of the system consists of two parts: (1) a sequence $(\alpha_1, \dots, \alpha_n)$ of n distinct elements in \mathbb{F}_{2^m} and (2) a square-free degree- t polynomial $g \in \mathbb{F}_{2^m}[x]$ such that $g(\alpha_i) \neq 0$ for all i .

For our implementation, g is generated as a uniform random degree- t monic irreducible polynomial in $\mathbb{F}_{2^m}[x]$. To generate g , we first generate a random element α in the extension field $\mathbb{F}_{2^{mt}}$. The polynomial g is then defined as the minimal polynomial of α in $\mathbb{F}_{2^m}[x]$, if the degree is t . To find the minimal polynomial, we view $\mathbb{F}_{2^{mt}}$ as the vector space $(\mathbb{F}_{2^m})^t$ and try to find linear dependency between $1, \alpha, \alpha^2, \dots, \alpha^t$ using Gaussian elimination. A description of the algorithm can be found in, for example, [26, Section 17.2].

The benefit of this approach is that it is easy to make Gaussian elimination constant time: [6] already shows how this can be achieved in the case of bit matrices. Note that the algorithm can fail to find a degree- t irreducible polynomial when $\alpha \in \mathbb{F}_{2^{mt}}$ such that t' is a divisor of t . For our parameters $m = 13$ and $t = 128$, the probability of failure is only 2^{-832} .

Recall that we use $n = 2^m$. Let ϕ be a permutation function such that $\phi(\hat{\alpha}_1, \dots, \hat{\alpha}_{2^m}) = (\alpha_1, \dots, \alpha_{2^m})$, where $(\hat{\alpha}_1, \dots, \hat{\alpha}_{2^m})$ is the standard order of field elements introduced by the FFT (see Sect. 4). In our software, the permutation function is defined using the condition bits in the corresponding size- 2^m Beneš network. We comment that there are $(2m - 1)2^{m-1} = m2^m - 2^{m-1}$ condition bits in the Beneš network, while a list of 2^m field elements takes $m2^m$ bits. In other words, representing $(\alpha_1, \dots, \alpha_n)$ as condition bits actually saves the size of secret keys.

Instead of generating the sequence α_i and then figure out the condition bits, in [14] the condition bits are simply generated as random bits in the current implementation. This approach is convenient for implementation. However, as discussed in [6, Section 5], this distribution of the sequence of field elements is not uniform, which might potentially lead to vulnerability (even though there has not been any evidence that the attacker can exploit the bias).

In order to avoid the problem, we decide to first generate a uniform random sequence and then figure out the corresponding condition bits. Generating the sequence can be done by first generating the list of all field elements and then attaching to each field element a random 32-bit integer. We then sort the list according to the 32-bit integers to obtain the sequence of field elements. The 32-bit numbers are checked so that there is no repetition; in case repetition is found, a new list of 32-bit numbers will be generated.

For figuring out the condition bits, we follow the approach given in the paper by Lev et al. [18]. The paper reduces the problem of generating the condition bits for a Beneš network to a coloring problem for a graph. A main building block in the Lev–Pippenger–Valiant algorithm is sorting. We use a piece of optimized constant-time code by Bernstein for Batcher’s odd-even sorting network [3] as a subroutine of condition-bit generation.

6.4 Public-key generation

Let H be the bit matrix obtained by replacing each entry in the matrix

$$\begin{pmatrix} 1/g(\alpha_1) & 1/g(\alpha_2) & \dots & 1/g(\alpha_n) \\ \alpha_1/g(\alpha_1) & \alpha_2/g(\alpha_2) & \dots & \alpha_n/g(\alpha_n) \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1^{t-1}/g(\alpha_1) & \alpha_2^{t-1}/g(\alpha_2) & \dots & \alpha_n^{t-1}/g(\alpha_n) \end{pmatrix}$$

by a column of m bits from the standard-basis representation. The receiver computes the row-reduced echelon form of H . If the result is of the form $[I|H']$, the public key is set to H' ; otherwise, a new secret key is generated.

In our implementation, the images $g(\hat{\alpha}_1), \dots, g(\hat{\alpha}_n)$ are first generated using the FFT implementation described in Sect. 4. After this, the inversions of all these images are computed, using Montgomery’s trick [21] with bitsliced field multiplications. Now we have the bitsliced representation of the first row of the matrix

$$\begin{pmatrix} 1/g(\hat{\alpha}_1) & 1/g(\hat{\alpha}_2) & \dots & 1/g(\hat{\alpha}_n) \\ \hat{\alpha}_1/g(\hat{\alpha}_1) & \hat{\alpha}_2/g(\hat{\alpha}_2) & \dots & \hat{\alpha}_n/g(\hat{\alpha}_n) \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\alpha}_1^{t-1}/g(\hat{\alpha}_1) & \hat{\alpha}_2^{t-1}/g(\hat{\alpha}_2) & \dots & \hat{\alpha}_n^{t-1}/g(\hat{\alpha}_n) \end{pmatrix}.$$

The remaining rows are then computed one by one using bitsliced field multiplications. Since all the rows are represented in the bitsliced format, the matrix can be easily viewed as the corresponding $mt \times n$ bit matrix. Then, the Beneš network is applied to each row of the bit matrix to obtain H . Finally, we follow [6, Section 6] to perform a constant-time Gaussian elimination. The public key is then the row-major representation of H' (one can of course use a column-major representation instead).

6.5 Encapsulation

The encapsulation begins with generating the error vector e of weight t . This is carried out by first generating a sequence of t random m -bit values, which indicates the positions of the errors. The t values are then checked for repetition. If a repetition is found, we simply regenerate the t random m -bit

Algorithm 1 The decoding procedure

```

1: function DECODING( $\hat{s}, g, \phi$ )
2:    $r \leftarrow (\hat{s}_1, \dots, \hat{s}_{mt}, 0, 0, \dots, 0) \in \mathbb{F}_2^n$ 
3:    $r \leftarrow \phi^{-1}(r)$ 
4:    $\beta = (\beta_1, \dots, \beta_n) \leftarrow \text{FFT}(g)$  ▷  $\beta \in \mathbb{F}_{2^m}^n$ 
5:    $s \leftarrow \text{FFT\_tr}(r_1/\beta_1^2, r_2/\beta_2^2, \dots, r_n/\beta_n^2)$  ▷  $s \in \mathbb{F}_{2^m}^{2t}$ 
6:    $f \leftarrow \text{BM}(s)$  ▷  $f \in \mathbb{F}_{2^m}[x]$ 
7:    $r' \leftarrow \text{FFT}(f)$  ▷  $f \in \mathbb{F}_{2^m}^n$ 
8:   Compute  $e \in \mathbb{F}_2^n$  such that  $e_i = 1$  iff  $r'_i = 0$ 
9:    $s' \leftarrow \text{FFT\_tr}(e_1/\beta_1^2, e_2/\beta_2^2, \dots, e_n/\beta_n^2)$  ▷  $s' \in \mathbb{F}_{2^m}^{2t}$ 
10:  if  $|e| \neq t$  or  $s \neq s'$  then return  $\perp$ 
11:  else return  $\phi(e)$ 
12:  end if
13: end function
    
```

values; otherwise, we convert the indices into the error vector as a sequence of $n/8$ bytes.

To compute each bit of the syndrome, each 128-bit word in the corresponding row is first ANDed with the corresponding 128-bit word in the error vector. The 128-bit results are then XORed together to form one single 128-bit word. We make use of the `popcnt` instruction to compute the parity of the 128-bit word, and the syndrome bit is set to the parity. Finally, after processing all the rows of the public key, we deal with the identity matrix by XORing the first $mt/8$ bytes of the error vector into the syndrome.

6.6 Decapsulation

As explained in [6], decoding consists of 5 stages: the initial permutation, syndrome computation, key-equation solving, root finding, and the final permutation. This is why the “total” column in [6, Table 1] is essentially

“perm” $\times 2$ + “synd” + “key eq” + “root”.

The “all” column in Table 1, however, is computed as

“perm” $\times 2$ + “synd” $\times 2$ + “key eq” + “root” $\times 2$.

In other words, we count one extra “root” and one extra “synd.”

The reason we count “root” one more time is a matter of implementation choice. To perform syndrome computation, each of the 2^m input bits is required to be scaled by $1/g(\alpha)^2$, where α is the corresponding point for evaluation. Since $1/g(\alpha)^2$ depends only on g , [6] uses them as pre-computed values. This strategy saves time but enlarges the size of secret keys. We decide to save the size of secret keys and compute all $1/g(\alpha)^2$ on the fly, using “root” for computing $g(\alpha)$, Montgomery’s trick for simultaneous inversions [21] with bitsliced multiplications, and bitsliced squarings.

The reason we count “synd” one more time is for re-encryption. A decoding algorithm is only required to decode

when the input syndrome corresponds to an error vector of weight t . For CCA-secure KEM, however, we need additionally the ability to reject invalid inputs. We therefore check the weight of the error vector and perform “synd” again to compute the syndrome of the error vector. The decoding is considered successful only if the weight is exactly t and the syndrome matches the output of the first “synd” stage.

We summarize the whole decoding procedure in Algorithm 1. The algorithm takes as input the syndrome \hat{s} (ciphertext for Niederreiter), the Goppa polynomial g , and the secret permutation ϕ . The algorithm’s output is either a weight- t error vector e which results in \hat{s} , or \perp .

The algorithm starts with generating a vector r that has syndrome \hat{s} with respected to $[I|H']$, by simply appending \hat{s} with 0’s. The next step is to compute the syndrome s of r with respect to

$$\begin{pmatrix} 1/g^2(\alpha_1) & 1/g^2(\alpha_2) & \dots & 1/g^2(\alpha_n) \\ \alpha_1/g^2(\alpha_1) & \alpha_2/g^2(\alpha_2) & \dots & \alpha_n/g^2(\alpha_n) \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_1^{2t-1}/g^2(\alpha_1) & \alpha_2^{2t-1}/g^2(\alpha_2) & \dots & \alpha_n^{2t-1}/g^2(\alpha_n) \end{pmatrix}.$$

This is achieved by first replacing r by $\phi^{-1}(r)$ (“perm”), scaling each entry r_i by $1/g^2(\hat{\alpha}_i)$, and then multiplying the result by

$$M = \begin{pmatrix} 1 & 1 & \dots & 1 \\ \hat{\alpha}_1 & \hat{\alpha}_2 & \dots & \hat{\alpha}_n \\ \vdots & \vdots & \ddots & \vdots \\ \hat{\alpha}_1^{2t-1} & \hat{\alpha}_2^{2t-1} & \dots & \hat{\alpha}_n^{2t-1} \end{pmatrix}.$$

The computation of $g(\hat{\alpha}_i)$ is carried out by the additive FFT (“root”), which is denoted as `FFT` in Algorithm 1, and the multiplication by M is carried out by the transposed additive FFT (“synd”), which is denoted as `FFT_tr`. The syndrome s is then fed into the Berlekamp-Massey algorithm (“key eq”) to generate the error locator f . From f , the corresponding error vector e is computed using the additive FFT (“root”). As we discussed above, it is necessary to check whether e ’s syndrome s' is the same as s . As for computing s , each entry e_i is first scaled by $1/g^2(\hat{\alpha}_i)$, and then, the scaled vector is multiplied by M to obtain s' . Note that only one transposed additive FFT (“synd”) is required here, as there is no need to compute $g(\hat{\alpha}_i)$ again. Finally, the algorithm outputs \perp if the weight of e is not t or if $s \neq s'$; otherwise, the algorithm outputs $\phi(e)$ (“perm”).

Acknowledgements This work is partially supported by Japan Society for the Promotion of Science KAKENHI Grant (C) (JP15K00183) and (JP15K00189) and Japan Science and Technology Agency, CREST (JPMJCR1404) and Infrastructure Development for Promoting International S&T Cooperation and Project for Establishing a Nationwide

Practical Education Network for IT Human Resources Development, Education Network for Practical Information Technologies. Permanent ID of this document: 9939a869c0a8d29b399e906a550f4cef. Date: 2018.02.23.

References

- Aragon, N., Barreto, P. S. L. M., Bettaieb, S., Bidoux, L., Blazy, O., Deneuville, J. C., Gaborit, P., Gueron, S., Güneysu, T., Melchor, C. A., Misoczki, R., Persichetti, E., Sendrier, N., Tillich, J. P., Zémor, G.: BIKE: Bit Flipping Key Encapsulation (2017). <http://bikesuite.org/files/BIKE.pdf>. Accessed 23 Feb 2018
- Avanzi, R., Bos, J., Ducas, L., Kiltz, E., Lepoint, T., Lyubashevsky, V., Schanck, J. M., Schwabe, P., Seiler, G., Stehlé, D.: CRYSTALS-KYBER, Algorithm Specifications and Supporting Documentation (2017). <https://pq-crystals.org/kyber/data/kyber-specification.pdf>. Accessed 23 Feb 2018
- Batcher, K. E.: Sorting networks and their applications. In: [7], pp. 307–314 (1968). <http://www.cs.kent.edu/~batcher/conf.html>. Accessed 23 Feb 2018
- Bernstein, D. J.: qhasm software package (2007). <http://cr.yo.to/qhasm.html>. Accessed 23 Feb 2018
- Bernstein, D. J., Chou, T., Lange, T., von Maurich, I., Misoczki, R., Niederhagen, R., Persichetti, E., Peters, C., Schwabe, P., Sendrier, N., Szefer, J., Wang, W.: Classic McEliece: Conservative Code-Based Cryptography (2017). <https://classic.mceliece.org/nist/mceliece-20171129.pdf>. Accessed 23 Feb 2018
- Bernstein, D. J., Chou, T., Schwabe, P.: McBits: fast constant-time code-based cryptography. In CHES 2013 [9], pp. 250–272 (2013)
- Bernstein, D. J., Chou, T., Schwabe, P.: AFIPS Conference Proceedings, vol. 32: 1968 Spring Joint Computer Conference, Reston, Virginia, Thompson Book Company (1968). See [3]
- Beneš, V.E.: Mathematical Theory of Connecting Networks and Telephone Traffic. Academic Press, London (1965)
- Bertoni, G., Coron, J. S. (eds.): Cryptographic Hardware and Embedded Systems—CHES 2013–15th International Workshop, Santa Barbara, CA, USA, Aug. 20–23, Proceedings, Lecture Notes in Computer Science, vol. 8086. Springer, Berlin (2013). ISBN 978-3-642-40348-4. See [6]
- Biham, E. (ed.): Fast Software Encryption, 4th International Workshop, FSE '97, Haifa, Israel, Jan. 20–22, Proceedings, Lecture Notes in Computer Science, vol. 1267. Springer, Berlin (1997). ISBN 3-540-63247-6. See [11]
- Biham, E.: A fast new DES implementation in software. In: [10], pp. 260–272 (1997)
- Biswas, B., Sendrier, N.: McEliece cryptosystem implementation: theory and practice. In: [13], pp. 47–62 (2008)
- Buchmann, J., Ding, J. (eds.): Post-quantum Cryptography, Second International Workshop, PQCrypto 2008, Cincinnati, OH, USA, Oct. 17–19, Proceedings, Lecture Notes in Computer Science, vol. 5299. Springer, Berlin (2008). See [12]
- Chou, T.: McBits revisited. In: CHES 2017 [15], pp. 213–231 (2017)
- Fischer, W., Homma, N. (eds.): Cryptographic Hardware and Embedded Systems—CHES 2017–19th International Conference, Taipei, Taiwan, Sept. 25–28, Proceedings, Lecture Notes in Computer Science, vol. 10529. Springer, Berlin (2017). ISBN 978-3-319-66786-7. See [14]
- Fog, A.: Instruction Tables (2017). http://www.agner.org/optimize/instruction_tables.pdf. Accessed 23 Feb 2018
- Gao, S., Mateer, T.: Additive fast Fourier transforms over finite fields. IEEE Trans. Inf. Theory **56**, 6265–6272 (2010). <http://www.math.clemson.edu/~sgao/pub.html>. Accessed 23 Feb 2018
- Lev, G. F., Pippenger, N., Valiant, L. G.: A fast parallel algorithm for routing in permutation networks. IEEE Trans. Comput. **C-30**, 93–100 (1981). <http://www.math.clemson.edu/~sgao/pub.html>. Accessed 23 Feb 2018
- Massey, J.L.: Shift-register synthesis and BCH decoding. IEEE Trans. Inf. Theory **15**, 122–127 (1969)
- McEliece, R. J.: A public-key cryptosystem based on algebraic coding theory. JPL DSN Progress Report, pp. 114–116 (1978). http://ipnpr.jpl.nasa.gov/progress_report2/42-44/44N.PDF. Accessed 23 Feb 2018
- Montgomery, P. L.: Speeding the Pollard and elliptic curve methods of factorization. Math. Comput. **48**, 243–264 (1987). <http://www.jstor.org/stable/pdf/2007888.pdf>. Accessed 23 Feb 2018
- Niederreiter, H.: Knapsack-type cryptosystems and algebraic coding theory. Probl. Control Inf. Theory **15**, 159–166 (1986)
- NIST: Submission Requirements and Evaluation Criteria for the Post-Quantum Cryptography Standardization Process (2016). <http://csrc.nist.gov/groups/ST/post-quantum-crypto/documents/call-for-proposals-final-dec-2016.pdf>. Accessed 23 Feb 2018
- Peters, C.: Information-set decoding for linear codes over F_q . In: PQCrypto 2010 [25], pp. 81–94 (2010). <http://eprint.iacr.org/2009/589>. Accessed 23 Feb 2018
- Sendrier, N. (ed.): Post-quantum Cryptography, Third International Workshop, PQCrypto, Darmstadt, Germany, May 25–28, Lecture Notes in Computer Science, vol. 6061. Springer, Berlin (2010). See [24]
- Shoup, V. (ed.): A Computational Introduction to Number Theory and Algebra (Version 2), Cambridge University Press, Cambridge. ISBN 978-0-521-51644-0 (2015). <http://www.shoup.net/ntb/ntb-v2.pdf>. Accessed 23 Feb 2018
- Youzhi, X.: Implementation of Berlekamp–Massey algorithm without inversion. IEE Proc. I Commun. Speech Vis. **138**, 138–140 (1991)