



## Efficient Arithmetic on Koblitz Curves\*

JEROME A. SOLINAS

National Security Agency, Ft. Meade, MD 20755, USA, and Visitor, Centre for Applied Cryptographic Research, Univ. of Waterloo

**Abstract.** It has become increasingly common to implement discrete-logarithm based public-key protocols on elliptic curves over finite fields. The basic operation is *scalar multiplication*: taking a given integer multiple of a given point on the curve. The cost of the protocols depends on that of the elliptic scalar multiplication operation.

Koblitz introduced a family of curves which admit especially fast elliptic scalar multiplication. His algorithm was later modified by Meier and Staffelbach. We give an improved version of the algorithm which runs 50% faster than any previous version. It is based on a new kind of representation of an integer, analogous to certain kinds of binary expansions. We also outline further speedups using precomputation and storage.

**Keywords:** elliptic curves, exponentiation, public-key cryptography

### 1. Introduction

It has become increasingly common to implement discrete-logarithm based public-key protocols on elliptic curves over finite fields. More precisely, one works with the points on the curve, which can be added and subtracted. If we add the point  $P$  to itself  $n$  times, we denote the result by  $nP$ . The operation of computing  $nP$  from  $P$  is called *scalar multiplication* by  $n$ . Elliptic public-key protocols are based on scalar multiplication, and the cost of executing such protocols depends mostly on the complexity of the scalar multiplication operation.

Scalar multiplication on an elliptic curve is analogous to exponentiation in the multiplicative group of integers modulo a fixed integer  $m$ . Various techniques have been developed [4] to speed modular exponentiation using memory and precomputations. Such methods, for the most part, carry over to elliptic scalar multiplication.

There are also efficiency improvements available in the elliptic case that have no analogue in modular exponentiation. There are three kinds of these:

1. One can choose the curve, and the base field over which it is defined, so as to optimize the efficiency of elliptic scalar multiplication. Thus, for example, one might choose the field of integers modulo a Mersenne prime, since modular reduction is particularly efficient [9] in that case. This option is not available for, say, RSA systems, since the secret primes are chosen randomly in order to maintain the security of the system.
2. One can use the fact that subtraction of points on an elliptic curve is just as efficient as addition. (The analogous statement for integers (mod  $m$ ) is false, since modular division is more expensive than modular multiplication.) The efficient methods for modular exponentiation all involve a sequence of squarings and multiplications that is

\* This paper is an expanded and updated version of the paper appearing in the Proceedings of Crypto '97.

based on the binary expansion of the exponent. The analogous procedure for elliptic scalar multiplication uses a sequence of doublings and additions of points. If we allow subtractions of points as well, we can replace [15] the binary expansion of the coefficient  $n$  by a more efficient *signed binary expansion* (i.e. an expansion in powers of two with coefficients 0 and  $\pm 1$ ).

3. One can use *complex multiplication*. Every elliptic curve over a finite field<sup>1</sup> comes equipped with a set of operations which can be viewed as multiplication by complex algebraic integers (as opposed to ordinary integers). These operations can be carried out efficiently for certain families of elliptic curves. In these cases, they can be utilized in various ways [10] to increase the efficiency of elliptic scalar multiplication.

It is the purpose of this paper to present a new technique for elliptic scalar multiplication. This new algorithm incorporates elements from all three of the above categories. The new method is 50% faster than any method previously known for operating on a non-supersingular elliptic curve.

## 2. Field and Elliptic Operations in $\mathbb{F}_{2^m}$

We begin with a brief survey of the various operations we will need in the field  $\mathbb{F}_{2^m}$  and on elliptic curves over this field.

*Squaring.* We will assume that the field  $\mathbb{F}_{2^m}$  is represented in terms of a *normal basis*: a basis over  $\mathbb{F}_2$  of the form

$$\{\theta, \theta^2, \theta^{2^2}, \dots, \theta^{2^{m-1}}\}.$$

The advantage of this representation is that squaring a field element can be accomplished by a one-bit cyclic shift of the bit string representing the element. This property will be crucial in what follows. If  $m$  is not divisible by 8, then one can use Gaussian cyclotomic periods to construct easily [1] an efficient normal basis for  $\mathbb{F}_{2^m}$ . (Since our application will require  $m$  to be prime, we can always use the Gaussian method.)

Our emphasis in this paper will be the case in which the field arithmetic is implemented in hardware. Although the algorithms that follow will be efficient in software as well, the full advantage of our method occurs in hardware, where the bit shifts (and therefore field squarings) are virtually free.

*Addition and Multiplication.* We may neglect the cost of additions in  $\mathbb{F}_{2^m}$  since they involve only bitwise XORs. A multiplication (of distinct elements) takes about  $m$  times as long, just as in the case of integer arithmetic. The cost of an elliptic operation depends mostly on the number of field multiplications it uses.

*Inversion.* Multiplicative inversion in  $\mathbb{F}_{2^m}$  can be performed in

$$L(m-1) + W(m-1) - 2$$

field multiplications using the method of [7]. Here  $L(k)$  represents the length of the binary expansion of  $k$ , and  $W(k)$  the number of ones in the expansion. This fact may be a consideration when choosing the degree  $m$ . (Alternatively, one can use the Euclidean algorithm [2], but one must first convert from the normal basis representation to the more familiar polynomial basis form, and then back again after the inversion.)

*Elliptic Addition.* The standard equation for an elliptic curve over  $\mathbb{F}_{2^m}$  is the *Weierstrass equation*

$$E: y^2 + xy = x^3 + ax^2 + b \quad (1)$$

where  $b \neq 0$ . Public key protocols based on this curve work on the group consisting of the points  $(x, y)$  on this curve, along with the group identity  $\mathcal{O}$ . (The element  $\mathcal{O}$  is called the *point at infinity*, but it is most convenient to represent it<sup>2</sup> by  $(0, 0)$ .)

ROUTINE 2 (ELLIPTIC GROUP OPERATION)

*Input:*

Points  $P_0$  and  $P_1$  on  $E$

*Output:*

The sum  $P_2 := P_0 + P_1$

*Computation:*

```

If  $P_0 = \mathcal{O}$  then output  $P_2 \leftarrow P_1$  and stop
If  $P_1 = \mathcal{O}$  then output  $P_2 \leftarrow P_0$  and stop
If  $x_0 = x_1$ 
  then
    if  $y_0 + y_1 = x_1$ 
      then
        output  $\mathcal{O}$  and stop
      else
        set  $\lambda \leftarrow x_1 + y_1/x_1$ 
            $x_2 \leftarrow \lambda^2 + \lambda + a$ 
            $y_2 \leftarrow x_1^2 + (\lambda + 1)x_2$ 
    else
      set  $\lambda \leftarrow (y_0 + y_1)/(x_0 + x_1)$ 
          $x_2 \leftarrow \lambda^2 + \lambda + x_0 + x_1 + a$ 
          $y_2 \leftarrow (x_1 + x_2)\lambda + x_2 + y_1$ 
Output  $P_2 \leftarrow (x_2, y_2)$ 

```

To *subtract* the point  $P = (x, y)$ , one adds the point  $-P = (x, x + y)$ .

Except for the special cases involving  $\mathcal{O}$ , the above addition and subtraction operations each require 1 multiplicative inversion and 2 multiplications.<sup>3</sup> (As always, we disregard the cost of adding and squaring field elements.)

### 3. Elliptic Scalar Multiplication

We next discuss the common methods for performing scalar multiplication on an arbitrary elliptic curve. These results will not be necessary for the subject of this paper, but will serve to motivate the new algorithms, which are analogues of these methods.

#### 3.1. The Addition-Subtraction Method

The basic technique for elliptic scalar multiplication is the *addition-subtraction method*. It is based on the *nonadjacent form* (NAF) of the coefficient  $n$ : a signed binary expansion with the property that no two consecutive coefficients are nonzero. For example,

$$\text{NAF}(29) = \langle 1, 0, 0, -1, 0, 1 \rangle \quad (3)$$

since  $29 = 32 - 4 + 1$ .

Just as every positive integer has a unique binary expansion, it also has a unique NAF. Moreover,  $\text{NAF}(n)$  has the fewest nonzero coefficients of any signed binary expansion of  $n$  [4]. There are several ways to construct the NAF of  $n$  from its binary expansion. We present the one that most resembles the new algorithm we will present in §4.

The idea is to divide repeatedly by 2. Recall that one can derive the binary expansion of an integer by dividing by 2, storing off the remainder (0 or 1), and repeating the process with the quotient. To derive a NAF, one allows remainders of 0 or  $\pm 1$ . If the remainder is to be  $\pm 1$ , one chooses whichever makes the quotient even.

ROUTINE 4 (NAF)

```

Input:
    a positive integer  $n$ 

Output:
     $\text{NAF}(n)$ 

Computation:
    Set  $c \leftarrow n$ 
    Set  $\mathcal{S} \leftarrow \langle \rangle$ 
    While  $c > 0$ 
        If  $c$  odd
            then
                set  $u \leftarrow 2 - (c \bmod 4)$ 
                set  $c \leftarrow c - u$ 
            else
                set  $u \leftarrow 0$ 
        Prepend  $u$  to  $\mathcal{S}$ 
        Set  $c \leftarrow c/2$ 
    EndWhile
    Output  $\mathcal{S}$ 

```

$c$	$u$	$\mathcal{S}$
29		$\langle \rangle$
28	1	
14		$\langle 1 \rangle$
14	0	
7		$\langle 0, 1 \rangle$
8	-1	
4		$\langle -1, 0, 1 \rangle$
4	0	
2		$\langle 0, -1, 0, 1 \rangle$
2	0	
1		$\langle 0, 0, -1, 0, 1 \rangle$
0	1	
0		$\langle 1, 0, 0, -1, 0, 1 \rangle$

Figure 1. Computing a NAF.

For example, to derive (3), one applies (4) with  $n = 29$ . The results are shown in Fig. 1.

Note that, although we have phrased the algorithm in terms of integer arithmetic, it can be implemented in terms of bit operations on the binary expansion of  $n$ . No arithmetic operations are needed beyond integer addition by 1.

In the derivation of the ordinary binary expansion, the sequence  $c$  is decreasing, but that is not true in general in (4). As a result, the NAF of a number may be longer than its binary expansion. Fortunately, it can be at most one bit longer, because

$$2^\ell < 3n < 2^{\ell+1} \quad (5)$$

where  $\ell$  is the bit length of  $\text{NAF}(n)$  [15].

The routine (4) can be modified as follows to produce an algorithm for elliptic scalar multiplication.

ROUTINE 6 (ADDITION-SUBTRACTION METHOD)

*Input:*

*a positive integer  $n$*   
*an elliptic curve point  $P$*

*Output:*

*the point  $nP$*

*Computation:*

Set  $c \leftarrow n$   
 Set  $Q \leftarrow \mathcal{O}$ ,  $P_0 \leftarrow P$   
 While  $c > 0$   
   If  $c$  odd then

```

    set  $u \leftarrow 2 - (c \bmod 4)$ 
    set  $c \leftarrow c - u$ 
    if  $u = 1$  then set  $Q \leftarrow Q + P_0$ 
    if  $u = -1$  then set  $Q \leftarrow Q - P_0$ 
  Set  $c \leftarrow c/2$ 
  Set  $P_0 \leftarrow 2P_0$ 
EndWhile
Output  $Q$ 

```

This algorithm is a *right-to-left* method, since (4) builds up the NAF starting at the least significant bit and ending at the most significant.<sup>4</sup> It is possible to give a *left-to-right* addition-subtraction method, but it has the disadvantage that it requires the entire NAF to be computed first, thus requiring more storage space and memory calls.

The cost of the addition-subtraction method depends on the bit length  $\ell$  of  $\text{NAF}(n)$ , which we now estimate. It follows from the Hasse theorem [22] that the order of an elliptic curve over  $\mathbb{F}_{2^m}$  is

$$\#E(\mathbb{F}_{2^m}) = 2^m + O(2^{m/2}).$$

Most public-key protocols on elliptic curves use a base point of prime order  $r$ . Since all of the curves (1) have even order, then

$$r \leq 2^{m-1} + O(2^{m/2}).$$

We can assume that  $n < r$ ; indeed, by using the identity

$$n(x, y) = (r - n)(x, x + y),$$

we can assume that  $n < r/2$ . Thus  $\ell < m$ , so that (6) requires about  $m$  doubles at most.

The number of additions is one less than the Hamming weight of (*i.e.* number of nonzero coefficients in)  $\text{NAF}(n)$ . The average density of nonzero coefficients among NAF's of length  $\ell$  is

$$\frac{2^\ell (3\ell - 4) - (-1)^\ell (6\ell - 4)}{9(\ell - 1)(2^\ell - (-1)^\ell)}, \quad (7)$$

or approximately (and asymptotically)  $1/3$  [15]. It follows via (5) that the Hamming weight  $H$  of  $\text{NAF}(n)$  satisfies

$$H \approx \frac{1}{3} \log_2 n. \quad (8)$$

Therefore, the average cost of (6) is  $\sim m$  doubles and  $\sim m/3$  additions, for a total of  $\sim 4m/3$  elliptic operations. This compares favorably to the classical *binary method*, which uses the ordinary binary expansion in place of the NAF. For binary expansions, the average density is exactly  $1/2$  rather than the value (7); thus the binary method requires about 12% more elliptic operations than the addition-subtraction method.

### 3.2. Window Methods

The addition-subtraction method can be generalized to produce still more efficient algorithms provided extra memory is available and precomputation is permitted. We present the basic method, called the *width- $w$  window method*.<sup>5</sup>

Let  $w$  be an integer greater than 1. Then each positive integer has a unique *width- $w$  NAF*: an expression

$$n = \sum_{j=0}^{\ell-1} u_j 2^j$$

where:

- each nonzero  $u_j$  is odd and less than  $2^{w-1}$  in absolute value;
- among any  $w$  consecutive coefficients, at most one is nonzero.

The case  $w = 2$  is that of the ordinary NAF.

The width- $w$  NAF is written

$$\text{NAF}_w(n) = \langle u_{\ell-1}, \dots, u_0 \rangle.$$

It can be computed via the following generalization of (4).

ROUTINE 9 (WIDTH- $w$  NAF)

*Input:*  
*a positive integer  $n$*

*Output:*  
 $\text{NAF}_w(n)$

*Computation:*  
 Set  $c \leftarrow n$   
 Set  $\mathcal{S} \leftarrow \langle \rangle$   
 While  $c > 0$   
   If  $c$  odd  
     then set  $u \leftarrow c \bmod 2^w$   
        set  $c \leftarrow c - u$   
     else set  $u \leftarrow 0$   
   Prepend  $u$  to  $\mathcal{S}$   
   Set  $c \leftarrow c/2$   
 EndWhile  
 Output  $\mathcal{S}$

The notation

$$a := b \bmod c$$

means that  $a$  is the integer satisfying

$$a \equiv b \pmod{c}$$

and

$$-\frac{c}{2} \leq a < \frac{c}{2}.$$

The routine (9) is commonly described by saying that one slides a “window” of width  $w$  along the binary expansion from right to left, using the contents to output the next entry of  $\text{NAF}_w(n)$ .

Given the width- $w$  NAF, one can perform elliptic scalar multiplication by  $n$  via the following algorithm.

ROUTINE 10 (WIDTH- $w$  ADDITION-SUBTRACTION METHOD)

*Input:*

*a positive integer  $n$*   
*an elliptic curve point  $P$*

*Output:*

*the point  $nP$*

*Precomputation:*

Set  $P_0 \leftarrow P$   
 Set  $P_{2^{w-2}-1} \leftarrow 2P_0$   
 For  $i$  from 1 to  $2^{w-2} - 1$  do  
   Set  $P_i \leftarrow P_{i-1} + P_{2^{w-2}-1}$   
 Next  $i$

*Computation:*

Compute  $\text{NAF}_w(n) = \langle u_{\ell-1}, \dots, u_0 \rangle$  (via (9))  
 Set  $Q \leftarrow \mathcal{O}$   
 For  $j$  from  $\ell - 1$  down to 0 do  
   Set  $Q \leftarrow 2Q$   
   If  $u_j \neq 0$  then  
      $i \leftarrow (|u_j| - 1)/2$   
     if  $u_j > 0$   
       then set  $Q \leftarrow Q + P_i$   
     else set  $Q \leftarrow Q - P_i$   
 Output  $Q$

The routine (10) is a left-to-right algorithm. The right-to-left algorithm (6) does not generalize well to the width- $w$  case, since each point  $P_i$  would have to be doubled  $\ell$  times. As remarked in [19], this is a general difficulty with window methods: the binary expansions must be computed right to left in general, but the elliptic scalar multiplication is best done

from left to right. This difficulty will not arise with the particular curves that are the subject of this paper.

The average density of a width- $w$  NAF is  $(w + 1)^{-1}$ . Thus one can diminish greatly the number of elliptic additions in an elliptic scalar multiplication, provided the memory is available. These speedups, however, do nothing to reduce the number of elliptic doublings. There seems to be little that can be done about this in general. The following special curves, however, admit elliptic scalar multiplication that do not use doublings at all.

#### 4. Koblitz Curves

The *anomalous binary curves* (or ABC's) are the curves  $E_0$  and  $E_1$  defined over  $\mathbb{F}_2$  by

$$E_a: y^2 + xy = x^3 + ax^2 + 1.$$

They are more commonly called *Koblitz curves* since their efficient scalar multiplication properties were first presented in [10].

##### 4.1. Basic Properties

We survey the basic properties of Koblitz curves that we will need.

*Group Orders.* We denote by  $E_a(\mathbb{F}_{2^m})$  the group of  $\mathbb{F}_{2^m}$ -rational points on  $E_a$ . This is the group on which the public-key protocols are performed. The group should be chosen so that it is computationally difficult to compute discrete logarithms of its elements. Thus, for example, the order  $\#E_a(\mathbb{F}_{2^m})$  should be divisible by a large prime [17]. Ideally,  $\#E_a(\mathbb{F}_{2^m})$  should be a prime or the product of a prime and small integer. This can only happen when  $m$  is itself prime, for otherwise there are large divisors arising from subgroups  $E_a(\mathbb{F}_{2^d})$  where  $d$  divides  $m$ .

When  $m$  is prime, the only such divisor is that arising from  $d = 1$ . The Koblitz curves over  $\mathbb{F}_2$  are

$$\begin{aligned} E_1(\mathbb{F}_2) &= \{\mathcal{O}, (0, 1)\} \\ E_0(\mathbb{F}_2) &= \{\mathcal{O}, (0, 1), (1, 0), (1, 1)\}. \end{aligned}$$

Since  $E_a(\mathbb{F}_2)$  is a subgroup of  $E_a(\mathbb{F}_{2^m})$ , it follows that the order  $\#E_a(\mathbb{F}_{2^m})$  is always divisible by

$$f = \#E_a(\mathbb{F}_2) = \begin{cases} 2 & \text{for } a = 1 \\ 4 & \text{for } a = 0. \end{cases} \quad (11)$$

We define an integer to be *very nearly prime* if it is of the form  $N = f \cdot r$ , where  $f = 2$  or 4 and  $r > 2$  is prime. Although the orders  $\#E_a(\mathbb{F}_{2^m})$  are never prime for  $m > 1$ , they are frequently very nearly prime. The values of  $m \leq 512$  for which  $\#E_1(\mathbb{F}_{2^m})$  is twice a prime are

$$m = 3, 5, 7, 11, 17, 19, 23, 101, 107, 109, 113, 163, 283, 311, 331, 347, 359.$$

The values of  $m \leq 512$  for which  $\#E_0(\mathbb{F}_{2^m})$  is 4 times a prime are

$$m = 5, 7, 13, 19, 23, 41, 83, 97, 103, 107, 131, 233, 239, 277, 283, 349, 409.$$

(The calculation of the orders is quite simple; the technique is given below.) The curves with very nearly prime order are the ones of most cryptographic interest.

*The Main Subgroup.* Suppose that  $\#E_a(\mathbb{F}_{2^m}) = f \cdot r$  is very nearly prime. We define the *main subgroup* to be the subgroup of order  $r$ . It is common<sup>6</sup> to perform cryptographic operations in the main subgroup rather than the entire curve.

**PROPOSITION 12** *Suppose that  $\#E_a(\mathbb{F}_{2^m})$  is very nearly prime, and let  $P$  be a point on  $E_a(\mathbb{F}_{2^m})$ . Then  $P$  is in the main subgroup if and only if  $P = fQ$  for some  $Q$  on  $E_a(\mathbb{F}_{2^m})$ .*

*Proof.* Both curves  $E_a(\mathbb{F}_2)$  are cyclic groups; to see this, one need only observe that  $2(1, 0) = (1, 1)$  on  $E_1(\mathbb{F}_2)$ . Thus the curve  $E_a(\mathbb{F}_{2^m})$  is cyclic whenever its order is very nearly prime. The result follows from the standard properties of finite cyclic groups. ■

As a result of (12), we have the following simple conditions to determine whether a given point is in the main subgroup. If  $a = 1$ , then a point  $P = (x, y)$  is in the main subgroup if and only if  $\text{Tr}(x) = 1$  [21]. If  $a = 0$ , then  $(x, y)$  is in the main subgroup if and only if  $\text{Tr}(x) = 0$  and  $\text{Tr}(y) = \text{Tr}(\lambda x)$ , where  $\lambda$  is an element with  $\lambda^2 + \lambda = x$ . (See Appendix A for proofs.) With a normal basis representation of the field, both the trace and the computation of  $\lambda$  can be done very efficiently. Thus, checking for membership in the main subgroup is essentially free if  $a = 1$ , and costs only one field multiplication if  $a = 0$ .

Because the main subgroup is the object of cryptographic interest, it is most important to optimize the elliptic scalar multiplication operation there. In §7.2, an algorithm for elliptic scalar multiplication in the main subgroup will be presented.

*Complex Multiplication.* Since the Koblitz curves are defined over  $\mathbb{F}_2$ , they have the following property: *if  $P = (x, y)$  is a point on  $E_a$ , then so is the point  $(x^2, y^2)$ .* Moreover, one can verify from (2) that

$$(x^4, y^4) + 2(x, y) = \mu \cdot (x^2, y^2) \tag{13}$$

for every  $(x, y)$  on  $E_a$ , where

$$\mu := (-1)^{1-a}. \tag{14}$$

This relation can be written more easily in terms of the Frobenius (squaring) map over  $\mathbb{F}_2$ :

$$\tau(x, y) := (x^2, y^2). \tag{15}$$

Using this notation, (13) becomes

$$\tau(\tau P) + 2P = \mu \tau P$$

for all  $P \in E_a$ . Symbolically, this can be written

$$(\tau^2 + 2)P = \mu \tau P.$$

This means that the squaring map can be regarded as implementing multiplication by the complex number  $\tau$  satisfying

$$\tau^2 + 2 = \mu \tau.$$

Explicitly, this number is

$$\tau = \frac{\mu + \sqrt{-7}}{2}.$$

By combining the squaring map with ordinary scalar multiplication, we can multiply points on  $E_a$  by any element of the ring  $\mathbb{Z}[\tau]$ . We say that  $E_a$  has *complex multiplication* by  $\tau$  [10].

*Lucas Sequences.* The *Lucas sequences* are sequences of integers that facilitate computations involving quadratic irrationals. The Lucas sequences for the numbers  $\tau$  will be used frequently in what follows; thus we summarize here their relevant properties.

- There are two Lucas sequences,  $U_k$  and  $V_k$ , associated with a given quadratic irrational. For  $\tau$ , the sequences are defined as follows.

$$\begin{aligned} U_0 = 0, \quad U_1 = 1 \quad \text{and} \quad U_{k+1} = \mu U_k - 2 U_{k-1} \quad \text{for } k \geq 1; \\ V_0 = 2, \quad V_1 = \mu \quad \text{and} \quad V_{k+1} = \mu V_k - 2 V_{k-1} \quad \text{for } k \geq 1. \end{aligned} \quad (16)$$

- It can be proved by induction that

$$\begin{aligned} U_k &= (\tau^k - \bar{\tau}^k) / \sqrt{-7} \\ V_k &= \tau^k + \bar{\tau}^k. \end{aligned} \quad (17)$$

If  $\theta := \tan^{-1}(\sqrt{7})$ , then (17) can be written as

$$\begin{aligned} U_k &= \mu^{k+1} 2^{k/2+1} \sin(k \theta) / \sqrt{7} \\ V_k &= \mu^k 2^{k/2+1} \cos(k \theta). \end{aligned} \quad (18)$$

- It can be proved by induction that

$$\tau^k = U_k \tau - 2 U_{k-1} \quad \text{for } k \geq 1. \quad (19)$$

Multiplying this equation by its conjugate, one obtains

$$U_k^2 - \mu U_k U_{k-1} + 2 U_{k-1}^2 = 2^{k-1} \quad \text{for } k \geq 1. \quad (20)$$

- The group orders  $\#E_a(\mathbb{F}_{2^m})$  are easily computed via

$$\#E_a(\mathbb{F}_{2^m}) = 2^m + 1 - V_m. \quad (21)$$

This identity follows from the basic properties of zeta functions of curves; [11].

*The Norm.* The norm of an element  $\alpha \in \mathbb{Z}[\tau]$  is the product of  $\alpha$  and its complex conjugate  $\bar{\alpha}$ . Explicitly, the norm of  $\delta := d_0 + d_1 \tau$  is

$$N(\delta) = d_0^2 + \mu d_0 d_1 + 2 d_1^2.$$

We will require the following properties of the norm.

- The norm function satisfies

$$N(\alpha \beta) = N(\alpha) N(\beta) \quad (22)$$

for all  $\alpha, \beta$  in  $\mathbb{Z}[\tau]$ .

- The Euclidean distance from  $\alpha$  to 0 in the complex plane is given by  $\sqrt{N(\alpha)}$ . Thus the Triangle Inequality takes the form

$$\sqrt{N(\alpha + \beta)} \leq \sqrt{N(\alpha)} + \sqrt{N(\beta)}. \quad (23)$$

- We will require the norms of some specific elements. It is easily checked that

$$N(\tau) = 2$$

and that

$$N(\tau - 1) = f$$

where  $f$  is as given in (11).

Finally,

$$N(\tau^m - 1) = \#E_a(\mathbb{F}_{2^m}) \quad (24)$$

(see [22]). (As an alternative to (21), the order of a Koblitz curve can be computed using this identity and (19).)

Since  $\tau - 1$  divides  $\tau^m - 1$ , it follows that  $f$  divides  $\#E_a(\mathbb{F}_{2^m})$  and that

$$N((\tau^m - 1)/(\tau - 1)) = \#E_a(\mathbb{F}_{2^m})/f. \quad (25)$$

- The ring  $\mathbb{Z}[\tau]$  is Euclidean with respect to the norm function [23]. That is, given an element  $\gamma$  and a nonzero element  $\delta$ , there exist elements  $\kappa$  and  $\rho$  such that  $\gamma = \delta \kappa + \rho$  and

$$N(\rho) < N(\delta). \quad (26)$$

As a result of this property, the ring  $\mathbb{Z}[\tau]$  has unique factorization. The element  $\tau$ , having prime norm, is a prime element.

#### 4.2. The $\tau$ -adic NAF

The complex multiplication property is useful for elliptic scalar multiplication because multiplication by  $\tau$ , being implemented by squaring, is essentially free when  $\mathbb{F}_{2^m}$  is represented in terms of a normal basis.<sup>7</sup> Thus it is worthwhile, when computing  $nP$ , to regard  $n$  as an element of  $\mathbb{Z}[\tau]$  rather than as “just” an integer. More precisely, one replaces the (signed) binary expansion of the coefficient with a (signed)  $\tau$ -adic expansion. That is, one represents  $n$  as a sum and difference of distinct powers of  $\tau$ .

For example, with  $a = 1$  we have

$$9 = \tau^5 - \tau^3 + 1. \quad (27)$$

Thus, if  $P = (x, y)$  is a point on  $E_1$ , then

$$9P = (x^{32}, y^{32}) - (x^8, y^8) + (x, y).$$

The above example gives 9 as what we call a  $\tau$ -adic NAF, since no two consecutive terms are nonzero. (Both [10] and [18] use signed  $\tau$ -adic expansions, but neither kind has the nonadjacency property.) As we shall see, the use of  $\tau$ -adic NAF’s gives a significant reduction in the number of terms, just as NAF’s give a significant improvement over binary expansions in the case of integers.

The remainder of this section is dedicated to proving the following result, which is analogous to the case of the ordinary NAF for integers.

**THEOREM 1** *Every element of the ring  $\mathbb{Z}[\tau]$  has a unique  $\tau$ -adic NAF.*

We can therefore speak of *the*  $\tau$ -adic NAF of an element  $\alpha$ . We will denote it by  $\text{TNAF}(\alpha)$ . Thus (27) is written

$$\text{TNAF}(9) = \langle 1, 0, -1, 0, 0, 1 \rangle.$$

In the course of proving Thm. 1, we will develop an efficient algorithm for computing the  $\tau$ -adic NAF for any element of  $\mathbb{Z}[\tau]$ .

**LEMMA 28** *The element  $c_0 + c_1 \tau$  of  $\mathbb{Z}[\tau]$  is divisible by  $\tau$  if and only if  $c_0$  is even. It is divisible by  $\tau^2$  if and only if*

$$c_0 \equiv 2c_1 \pmod{4}. \quad (29)$$

*Proof.* If  $\mu$  is defined as in (14), then

$$(d_0 + d_1 \tau) \tau = -2d_1 + (d_0 + \mu d_1) \tau.$$

It follows at once that, if  $\tau$  divides  $c_0 + c_1 \tau$ , then  $c_0$  is even.

Conversely, if  $c_0$  is even, then

$$\frac{c_0 + c_1 \tau}{\tau} = \frac{\mu c_0 + 2 c_1}{2} - \frac{c_0}{2} \tau$$

is an element of  $\mathbb{Z}[\tau]$ .

To prove the corresponding statements for  $\tau^2$ , we begin with the identity

$$\tau^2 = \mu \tau - 2.$$

It follows that every multiple of  $\tau^2$  has the form

$$(d_0 + d_1 \tau) (\mu \tau - 2) = -2 (d_0 + \mu d_1) + (\mu d_0 - d_1) \tau.$$

It is easily verified that the values

$$c_0 = -2 (d_0 + \mu d_1)$$

$$c_1 = \mu d_0 - d_1$$

satisfy (29). Conversely,

$$\frac{c_0 + c_1 \tau}{\tau^2} = -\frac{(1 + 2\mu) c_0 + 2\mu c_1}{4} + \mu \cdot \frac{c_0 - 2 c_1}{4} \tau,$$

which is easily seen to be an element of  $\mathbb{Z}[\tau]$  if (29) holds. ■

In light of (28), we have the following algorithm (joint work with R. Reiter) for computing the  $\tau$ -adic NAF. It is completely analogous to (4), but here we are dividing by  $\tau$  rather than by 2. Since  $\tau$  has norm 2, then by (26), the possible remainders upon division by  $\tau$  are  $\pm 1$ . Earlier algorithms chose the remainder that minimized the norm of the quotient; this is analogous to the basic division algorithm for generating the binary expansion of an integer. What we shall do instead is to choose the remainder that makes the quotient divisible by  $\tau$ . This is analogous to the computation of the NAF for integers.

ALGORITHM 1 ( $\tau$ -adic NAF)

*Input:*

integers  $r_0, r_1$

*Output:*

TNAF( $r_0 + r_1 \tau$ )

*Computation:*

Set  $c_0 \leftarrow r_0, c_1 \leftarrow r_1$

Set  $\mathcal{S} \leftarrow \langle \rangle$

While  $c_0 \neq 0$  or  $c_1 \neq 0$

If  $c_0$  odd

then

$c_0$	$c_1$	$u$	$\mathcal{S}$
9	0		$\langle \rangle$
8	0	1	
4	-4		$\langle 1 \rangle$
4	-4	0	
-2	-2		$\langle 0, 1 \rangle$
-2	-2	0	
-3	1		$\langle 0, 0, 1 \rangle$
-2	1	-1	
0	1		$\langle -1, 0, 0, 1 \rangle$
0	1	0	
1	0		$\langle 0, -1, 0, 0, 1 \rangle$
0	0	1	
0	0		$\langle 1, 0, -1, 0, 0, 1 \rangle$

Figure 2. Computing a  $\tau$ -adic NAF.

```

        set  $u \leftarrow 2 - (c_0 - 2c_1 \bmod 4)$ 
        set  $c_0 \leftarrow c_0 - u$ 
    else
        set  $u \leftarrow 0$ 
    Prepend  $u$  to  $\mathcal{S}$ 
    Set  $(c_0, c_1) \leftarrow (c_1 + \mu c_0/2, -c_0/2)$ 
EndWhile
Output  $\mathcal{S}$ 

```

For example, to derive (27), one applies Alg. 1 with  $a = 1$ ,  $c_0 = 9$ , and  $c_1 = 0$ . The results are shown in Fig. 2.

Note that the implementation of Alg. 1 involves nothing more complicated than integer addition. (This is slightly more than is required by (4), which only adds 1 to an integer.)

Having described how to compute a  $\tau$ -adic NAF, we now prove that such a representation is unique.

LEMMA 30 *Let  $\alpha \in \mathbb{Z}[\tau]$ . Then precisely one of the following statements hold:*

- $\alpha$  is divisible by  $\tau$ .
- $\alpha \equiv 1 \pmod{\tau^2}$ .
- $\alpha \equiv -1 \pmod{\tau^2}$ .

*Proof.* Let  $\alpha = c_0 + c_1 \tau$ . If  $c_0$  is even, then  $\tau$  divides  $\alpha$  by (28). If  $c_0$  is odd, then precisely one of  $c_0 \pm 1$  satisfies (29); thus  $\tau^2$  divides precisely one of  $\alpha \pm 1$ , by (28). ■

**COROLLARY 31** *Let  $\alpha \in \mathbb{Z}[\tau]$ . Then any two  $\tau$ -adic NAF's of  $\alpha$  have the same rightmost entry.*

*Proof.* Depending on which of the three possibilities hold for  $\alpha$ , the rightmost entry of a  $\tau$ -adic NAF for  $\alpha$  is 0, 1, or  $-1$ , respectively. ■

*Proof of Thm. 1.* The existence of the  $\tau$ -adic NAF is established by Alg. 1. The uniqueness follows from (31), using induction on the length of the  $\tau$ -adic NAF. ■

### 4.3. Length and Density of $\tau$ -adic NAF's

To evaluate the usefulness of the  $\tau$ -adic NAF as a substitute for an ordinary NAF, it is necessary to know its Hamming weight (*i.e.* number of nonzero terms). In the case of the ordinary NAF, the weight is calculated in terms of the length of the NAF (given by (5)) and the density of nonzero terms (given by (7)). We now obtain analogous results for  $\tau$ -adic NAF's.

The calculation of the density is trivial.

**PROPOSITION 32** *The average density among  $\tau$ -adic NAF's of length  $\ell$  is given by (7), and is therefore asymptotically  $1/3$ .*

*Proof.* The result follows from (7) since the same sequences occur as length- $\ell$  NAF's and length- $\ell$   $\tau$ -adic NAF's. ■

The estimation of the length of a NAF is more involved. Intuitively, the answer should be

$$\ell \approx \log_2(N(\alpha)),$$

because Alg. 1 begins with  $\alpha$ , and divides by the norm-2 element  $\tau$  in computing each entry of the  $\tau$ -adic NAF. Of course this ignores the effect of the additions and subtractions, but provides a benchmark for the results to follow.

We begin with some notation and terminology. By the *length* of an element of  $\mathbb{Z}[\tau]$  we mean the length of its  $\tau$ -adic NAF. Let  $N_{\max}(k)$  denote the largest norm occurring among all length- $k$  elements of  $\mathbb{Z}[\tau]$ .

**LEMMA 33** *For all  $k$ ,*

$$2 N_{\max}(k) \leq N_{\max}(k + 1).$$

*Proof.* If  $\alpha$  is a length- $k$  element of maximal norm, then  $\tau\alpha$  is an element of length  $k + 1$  and norm  $2 N(\alpha)$ . ■

**COROLLARY 34**  *$N_{\max}(k)$  is the largest norm occurring among all elements  $\alpha$  of  $\mathbb{Z}[\tau]$  whose length is at most  $k$ .*

LEMMA 35 *If  $c > e$ , then*

$$\sqrt{N_{\max}(c)} \leq 2^{e/2} \sqrt{N_{\max}(c-e)} + \sqrt{N_{\max}(e)}.$$

*Proof.* Let  $\gamma$  be a length- $c$  element of maximal norm, so that

$$N(\gamma) = N_{\max}(c). \quad (36)$$

Let the  $\tau$ -adic NAF of  $\gamma$  be

$$\text{TNAF}(\gamma) = \langle u_{c-1}, \dots, u_0 \rangle, \quad (37)$$

and define  $\rho$  by

$$\text{TNAF}(\rho) = \langle u_{e-1}, \dots, u_0 \rangle. \quad (38)$$

Then  $\rho$  has length at most  $e$ , so that

$$N(\rho) \leq N_{\max}(e) \quad (39)$$

by (34). Finally,

$$\gamma = \tau^e \delta + \rho \quad (40)$$

for some  $\delta$  of length  $c - e$ ; note that

$$N(\delta) \leq N_{\max}(c - e). \quad (41)$$

It follows from the Triangle Inequality (23) that

$$\sqrt{N(\gamma)} \leq 2^{e/2} \sqrt{N(\delta)} + \sqrt{N(\rho)}.$$

The result now follows by (36), (39), and (41). ■

For convenience we introduce the notation

$$M_k := \sqrt{2^{-k} N_{\max}(k)}.$$

In this notation, (33) states that

$$M_k \leq M_{k+1} \quad (42)$$

for all  $k$ , and (35) states that

$$M_c \leq M_{c-e} + 2^{-(c-e)/2} M_e \quad (43)$$

if  $c > e$ .

PROPOSITION 44 For  $d$  and  $q$  positive,

$$\frac{M_{dq}}{1 - 2^{-dq/2}} \leq \frac{M_d}{1 - 2^{-d/2}}.$$

*Proof.* It follows from (43) that

$$M_{(k+1)d} - M_{kd} \leq 2^{-kd/2} M_d$$

for  $k = 1, \dots, (q - 1)$ . Summing over  $k$ , we obtain

$$M_{dq} - M_d \leq M_d \sum_{k=1}^{q-1} (2^{-d/2})^k.$$

The result follows from summing the geometric series. ■

PROPOSITION 45 For  $\ell > d$ ,

$$M_\ell < \frac{M_d}{1 - 2^{-d/2}}.$$

*Proof.* Let  $dq$  be the smallest multiple of  $d$  greater than or equal to  $\ell$ . Then  $M_\ell \leq M_{dq}$  by (42), and

$$M_{dq} < \frac{M_d}{1 - 2^{-d/2}}$$

by (44). ■

COROLLARY 46 For  $\ell > d$ ,

$$N_{\max}(\ell) < \frac{N_{\max}(d)}{(2^{d/2} - 1)^2} \cdot 2^\ell.$$

The result (46) provides the basic upper bound for the norm of a length- $\ell$  element. We now derive a lower bound. Let  $N_{\min}(k)$  denote the smallest norm occurring among all length- $k$  elements of  $\mathbb{Z}[\tau]$ .

LEMMA 47 If  $c > e$ , then

$$\sqrt{N_{\min}(c)} \geq 2^{e/2} \sqrt{N_{\min}(c - e)} - \sqrt{N_{\max}(e)}.$$

*Proof.* Let  $\gamma$  be a length- $c$  element of minimal norm, so that

$$N(\gamma) = N_{\min}(c). \tag{48}$$

Let the  $\tau$ -adic NAF of  $\gamma$  be as in (37). Define  $\rho$  as in (38), and  $\delta$  as in (40). Then

$$N(\delta) \geq N_{\min}(c - e) \tag{49}$$

and

$$N(\rho) \leq N_{\max}(e), \quad (50)$$

the latter following from (34). By the Triangle Inequality (23),

$$\sqrt{N(\gamma)} \geq 2^{e/2} \sqrt{N(\delta)} - \sqrt{N(\rho)}.$$

The result now follows by (48), (49), and (50).  $\blacksquare$

**COROLLARY 51** For  $\ell > 2d$ ,

$$N_{\min}(\ell) > \left( \sqrt{N_{\min}(d)} - \frac{\sqrt{N_{\max}(d)}}{2^{d/2} - 1} \right)^2 \cdot 2^{\ell-d}.$$

*Proof.* Follows from (46) and (47).  $\blacksquare$

Combining (46) and (51), we obtain the main result of this section.

**THEOREM 2** Let  $\ell > 2d$ , and let  $\alpha$  be a length- $\ell$  element of  $\mathbb{Z}[\tau]$ . Then

$$\left( \sqrt{N_{\min}(d)} - \frac{\sqrt{N_{\max}(d)}}{2^{d/2} - 1} \right)^2 \cdot 2^{\ell-d} < N(\alpha) < \frac{N_{\max}(d)}{(2^{d/2} - 1)^2} \cdot 2^\ell.$$

To apply Thm. 2, we choose a small value of  $d$ , and evaluate  $N_{\max}(d)$  and  $N_{\min}(d)$  by direct evaluation of all length- $d$  elements of  $\mathbb{Z}[\tau]$ . This is only feasible if  $d$  is quite small, but the resulting bounds are quite accurate.

For example, we apply Thm. 2 with  $d = 15$ . The bounds are

$$N_{\max}(15) = 47324 \quad \text{and} \quad N_{\min}(15) = 2996.$$

It follows that

$$1.399009614 \cdot 2^{\ell-4} < N(\alpha) < .7301517653 \cdot 2^{\ell+1}. \quad (52)$$

In other words, the length of the  $\tau$ -adic NAF is bounded by

$$\log_2(N(\alpha)) - .5462682713 < \ell < \log_2(N(\alpha)) + 3.51559412 \quad (53)$$

when  $\ell > 30$ .

To measure the quality of these bounds, we compare them with the norms and lengths of some specific (infinite) families of elements.

Let  $a = 1$  and

$$\beta := \tau^{14} + \tau^{11} - \tau^9 + \tau^7 - \tau^4 + \tau.$$

Then  $N(\beta) = 2996$ , minimal among elements of length 15. Since

$$N(\beta) \approx (1.4628906) \cdot 2^{11},$$

it follows that, for any  $k \geq 0$ , the element  $\alpha := \beta \tau^k$  of length  $\ell := k + 15$  satisfies

$$N(\alpha) \approx (1.4628906) \cdot 2^{\ell-4}.$$

This exceeds the lower bound given in (52) by only a small amount.

Now let  $a = 1$  and

$$\beta := \tau^{11} - \tau^9 + \tau^6 - \tau^4 + \tau.$$

Then  $N(\beta) = 5842$ , maximal among elements of length up to 12. Let

$$\gamma := (\tau^{13} - 1) \beta.$$

Clearly, the  $\tau$ -adic NAF of  $\gamma$  is the concatenation of two copies of that of  $\beta$ , separated by a zero. Thus  $\gamma$  has length 25. Its norm is

$$N(\gamma) = 5842 \cdot N(\tau^{13} - 1)$$

by (22). Now

$$N(\tau^{13} - 1) = 2^{13} + 1 - V_{13}$$

by (24) and (21). Computing the Lucas element  $V_{13}$  via (16), we find that

$$N(\tau^{13} - 1) = 8374,$$

so that

$$N(\gamma) \approx (.7289783) \cdot 2^{26}.$$

It follows that, for any  $k \geq 0$ , the element  $\alpha := \gamma \tau^k$  of length  $\ell := k + 25$  satisfies

$$N(\alpha) \approx (.7289783) \cdot 2^{\ell+1}.$$

This is just slightly under the upper bound given in (52).

It follows from (32) that the average Hamming weight among length- $\ell$  TNAF's is roughly  $\ell/3$ . If we assume that this average value holds for the subset of TNAF's of rational integers, then it follows from (53) that the Hamming weight  $H$  of the  $\tau$ -adic NAF for the integer  $n$  satisfies

$$H \approx \frac{2}{3} \log_2 n. \tag{54}$$

This is twice as large as the Hamming weight of an ordinary NAF, because the  $\tau$ -adic NAF is twice as long. If we replace the ordinary NAF by the  $\tau$ -adic NAF, then, we will have eliminated the elliptic doublings in our scalar multiplication method, but doubled the number of elliptic additions. This largely mitigates the advantage of the  $\tau$ -adic method.

Fortunately, this situation can be fixed. The solution is to replace the  $\tau$ -adic NAF by an equivalent expression, called a *reduced  $\tau$ -adic NAF*, that is only half as long.

Before presenting this, however, it is necessary to develop the machinery of modular reduction in the ring  $\mathbb{Z}[\tau]$ .

### 5. Modular Reduction in $\mathbb{Z}[\tau]$

In this section, we define precisely what is meant by modular reduction in the ring  $\mathbb{Z}[\tau]$ , and present an efficient method for implementing it.

Our technique is a generalization of the notion of modular reduction in the ring  $\mathbb{Z}$  of rational integers. Suppose  $c$  and  $d > 1$  are integers. It is desired to reduce  $c$  modulo  $d$ , *i.e.* find the integer

$$\rho := c \text{ mods } d$$

where the “mods” notation is as in §3.2. The integer  $\rho$  can be found by integer division: if

$$\text{Round}(\lambda) := \lfloor \lambda + 1/2 \rfloor,$$

then

$$\rho = c - \kappa d,$$

where

$$\kappa := \text{Round}(c/d).$$

A more compact way of describing  $\rho$  is in terms of the *fractional part* operation. The fractional part of  $\lambda$  is defined to be

$$((\lambda)) := \lambda - \text{Round}(\lambda).$$

The modular reduction process can then be described by

$$\rho = d \left( \left( \frac{c}{d} \right) \right). \quad (55)$$

Since

$$\left( \left( \frac{c}{d} \right) \right) < \frac{1}{2},$$

it follows that

$$N(\rho) < \frac{N(d)}{2}.$$

We now generalize these concepts to the ring  $\mathbb{Z}[\tau]$ .

#### 5.1. Rounding and Fractional Parts in $\mathbb{Z}[\tau]$

We begin by extending to  $\mathbb{Z}[\tau]$  the definitions of  $\text{Round}(\lambda)$  and  $((\lambda))$ . (The variable  $\lambda$  now denotes the expression  $\lambda_0 + \lambda_1 \tau$ , where the  $\lambda_i$  are real numbers.)

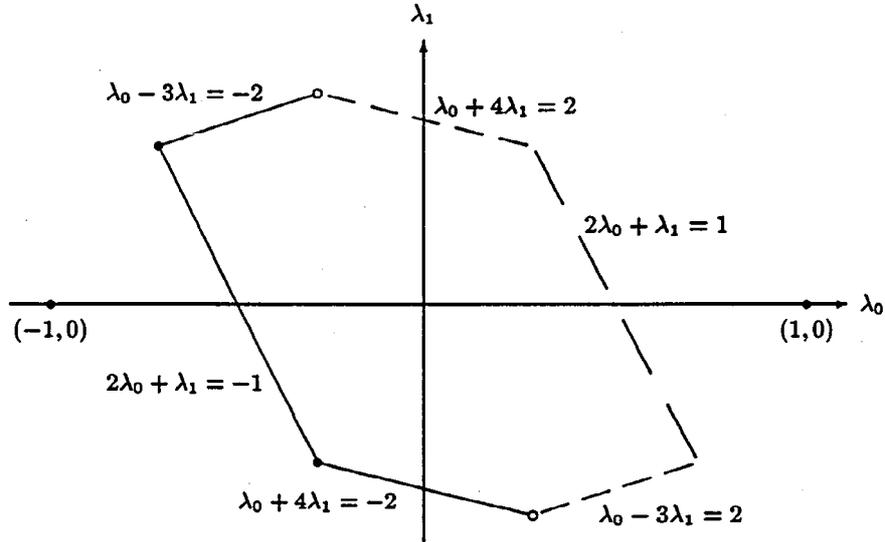


Figure 3. The region  $\mathcal{U}$  for the case  $a = 1$ .

We define  $\mathcal{U}$  to be the region in the  $(\lambda_0, \lambda_1)$ -plane given by the inequalities

$$\begin{aligned} -1 &\leq 2\lambda_0 + \mu\lambda_1 < 1 \\ -2 &\leq \lambda_0 + 4\mu\lambda_1 < 2 \\ -2 &\leq \lambda_0 - 3\mu\lambda_1 < 2. \end{aligned} \tag{56}$$

(See Fig. 3 for  $a = 1$  and Fig. 4 for  $a = 0$ .) Copies of  $\mathcal{U}$  tile the plane (see Fig. 5 for the case  $a = 1$ ), with each copy having as its center an element of  $\mathbb{Z}[\tau]$ . Given  $\lambda \in \mathbb{Q}(\tau)$ , we “round off” by choosing as  $\kappa \in \mathbb{Z}[\tau]$  the center of the copy of  $\mathcal{U}$  containing  $\lambda$ . We will denote this operation either by

$$(q_0, q_1) = \text{Round}(\lambda_0, \lambda_1)$$

or by

$$q_0 + q_1 \tau = \text{Round}(\lambda_0 + \lambda_1 \tau),$$

since there is no danger of confusion.

As in the integer case, we define

$$((\lambda)) := \lambda - \text{Round}(\lambda)$$

for all complex  $\lambda$ . It is easy to see that the set of possible values for  $((\lambda))$  is precisely the region  $\mathcal{U}$ .

We next prove the main properties of these operations.

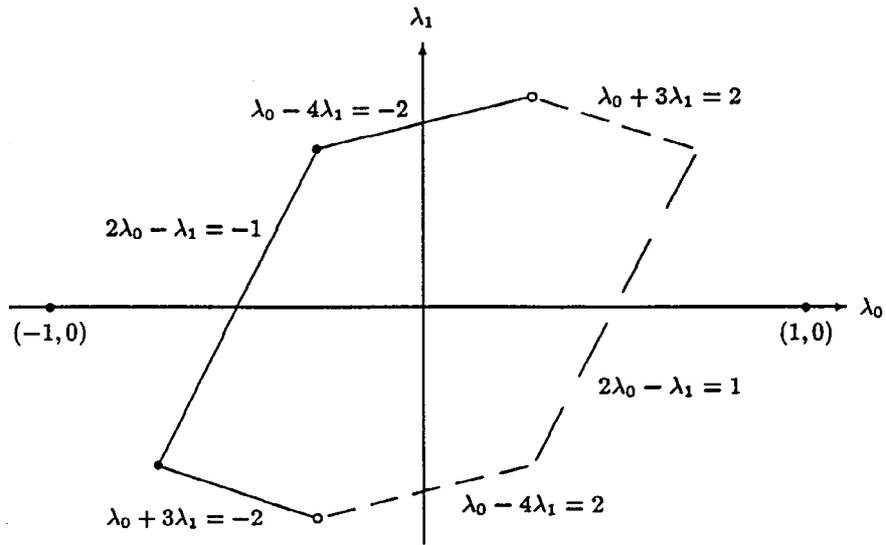


Figure 4. The region  $\mathcal{U}$  for the case  $a = 0$ .

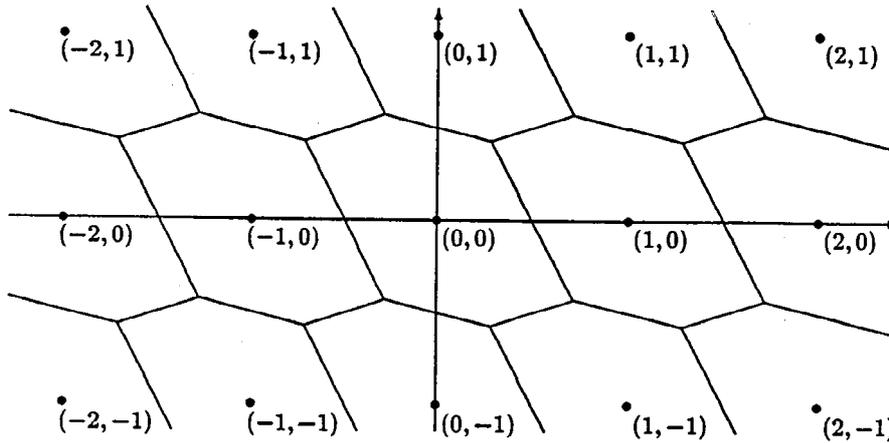


Figure 5. Copies of  $\mathcal{U}$  for the case  $a = 1$ .

PROPOSITION 57 *Suppose that  $\lambda$  is in the interior of the region  $\mathcal{U}$ . Then*

$$N(\lambda) < \frac{4}{7}.$$

*Proof.* The set of points in the  $(\lambda_0, \lambda_1)$ -plane of norm  $4/7$  forms the ellipse

$$\lambda_0^2 + \mu \lambda_0 \lambda_1 + 2 \lambda_1^2 = \frac{4}{7}.$$

Now each of the six vertices of  $\mathcal{U}$  has norm  $4/7$ , and so lies on the ellipse. Since  $\mathcal{U}$  is convex, it lies entirely in the ellipse. The result follows from the fact that the points inside the ellipse are those of norm less than  $4/7$ . ■

PROPOSITION 58 *Suppose that  $\lambda$  is in the interior of the region  $\mathcal{U}$ . Then*

$$N(\lambda) < N(\lambda + \alpha)$$

*for every nonzero  $\alpha \in \mathbb{Z}[\tau]$ .*

*Proof.* It is straightforward to check that

$$N(\lambda) < N(\lambda \pm 1) \quad \text{if and only if} \quad |2\lambda_0 + \mu\lambda_1| < 1;$$

that

$$N(\lambda) < N(\lambda \pm \tau) \quad \text{if and only if} \quad |\mu\lambda_0 + 4\lambda_1| < 2;$$

and that

$$N(\lambda) < N(\lambda \pm \bar{\tau}) \quad \text{if and only if} \quad |\mu\lambda_0 - 3\lambda_1| < 2.$$

Since  $\lambda$  lies in the interior of  $\mathcal{U}$ , then by (56) it satisfies all three conditions. Thus the result is proved for  $\alpha = \pm 1, \pm\tau, \pm\bar{\tau}$ . Now let  $\alpha$  be any other nonzero element of  $\mathbb{Z}[\tau]$ . Then  $N(\alpha) \geq 4$ ; and  $N(\lambda) < 4/7$  by (57). Thus the result in this case follows from (23). ■

The following properties follow from (57) and (58).

COROLLARY 59 *If  $\kappa := \text{Round}(\lambda)$  and  $\zeta := ((\lambda)) = \lambda - \kappa$ , then:*

- $N(\zeta) \leq N(\zeta + \alpha)$  for every  $\alpha \in \mathbb{Z}[\tau]$ .
- $N(\zeta) \leq \frac{4}{7}$ .

The first condition of (59) simply says that  $\text{Round}(\lambda)$  is the element of  $\mathbb{Z}[\tau]$  closest to  $\lambda$ , justifying the terminology. The second condition was proved in [18].

We now give an algorithm for computing  $\text{Round}(\lambda)$ .

## ROUTINE 60 (ROUNDING OFF)

*Input:*  
*real numbers  $\lambda_0, \lambda_1$  specifying the complex number  $\lambda := \lambda_0 + \lambda_1 \tau$*

*Output:*  
*real numbers  $q_0, q_1$  specifying  $q_0 + q_1 \tau := \text{Round}(\lambda)$*

```

Set  $f_0 \leftarrow \text{Round}(\lambda_0)$ 
     $f_1 \leftarrow \text{Round}(\lambda_1)$ 
Set  $\eta_0 \leftarrow \lambda_0 - f_0$ 
     $\eta_1 \leftarrow \lambda_1 - f_1$ 
Set  $h_0 \leftarrow 0$ 
     $h_1 \leftarrow 0$ 
Set  $\eta \leftarrow 2\eta_0 + \mu\eta_1$ 
If  $\eta \geq 1$ 
  then
    if  $\eta_0 - 3\mu\eta_1 < -1$ 
      then set  $h_1 \leftarrow \mu$ 
      else set  $h_0 \leftarrow 1$ 
    else
      if  $\eta_0 + 4\mu\eta_1 \geq 2$ 
        then set  $h_1 \leftarrow \mu$ 
  If  $\eta < -1$ 
  then
    if  $\eta_0 - 3\mu\eta_1 \geq 1$ 
      then set  $h_1 \leftarrow -\mu$ 
      else set  $h_0 \leftarrow -1$ 
    else
      if  $\eta_0 + 4\mu\eta_1 < -2$ 
        then set  $h_1 \leftarrow -\mu$ 
Set  $q_0 \leftarrow f_0 + h_0$ 
     $q_1 \leftarrow f_1 + h_1$ 
Output  $q_0, q_1$ 

```

**5.2. Division and Modular Reduction in  $\mathbb{Z}[\tau]$** 

We now use the rounding-off operation to develop algorithms for division and modular reduction in  $\mathbb{Z}[\tau]$ .

The analogue of integer division is as follows. Given a dividend  $\gamma = c_0 + c_1 \tau$  and a divisor  $\delta = d_0 + d_1 \tau$ , we wish to find a quotient  $\kappa = q_0 + q_1 \tau$  and a remainder  $\rho = r_0 + r_1 \tau$ , such that

$$\gamma = \kappa \delta + \rho$$

and such that  $\rho$  is as small (in norm) as possible. To do this, we obtain  $\kappa$  by rounding off  $\gamma/\delta$  and then solving for  $\rho$ . That is, we let

$$\begin{aligned}\lambda &:= \frac{\gamma}{\delta} \\ &= \frac{\gamma \bar{\delta}}{N(\delta)} \\ &= \frac{g_0 + g_1 \tau}{N},\end{aligned}$$

find  $\kappa$  via

$$\kappa := \text{Round}\left(\frac{g_0}{N} + \frac{g_1}{N} \tau\right), \quad (61)$$

and obtain  $\rho$  via

$$\rho := \gamma - \kappa \delta.$$

The following algorithm gives these steps explicitly in terms of the coefficients.

ROUTINE 62 (DIVISION IN  $\mathbb{Z}[\tau]$ )

*Input:*

*the dividend  $\gamma = c_0 + c_1 \tau$  and divisor  $\delta = d_0 + d_1 \tau$*

*Output:*

*the quotient  $\kappa = q_0 + q_1 \tau$  and the remainder  $\rho = r_0 + r_1 \tau$*

*Computation:*

$$\begin{aligned}g_0 &\leftarrow c_0 d_0 + \mu c_0 d_1 + 2 c_1 d_1 \\ g_1 &\leftarrow c_1 d_0 - c_0 d_1 \\ N &\leftarrow d_0^2 + \mu d_0 d_1 + 2 d_1^2 \\ \lambda_0 &\leftarrow g_0/N \\ \lambda_1 &\leftarrow g_1/N \\ (q_0, q_1) &\leftarrow \text{Round}(\lambda_0, \lambda_1) \\ r_0 &\leftarrow c_0 - d_0 q_0 + 2 d_1 q_1 \\ r_1 &\leftarrow c_1 - d_1 q_0 - d_0 q_1 - \mu d_1 q_1 \\ \text{Output } &q_0, q_1, r_0, r_1\end{aligned}$$

If we disregard the quotient  $\kappa$  and only output the remainder  $\rho$ , this routine may be viewed as a modular reduction algorithm. In this case, we write

$$\rho := \gamma \bmod \delta.$$

In analogy with (55) we have

$$\rho = \delta \left( \left( \frac{\gamma}{\delta} \right) \right). \quad (63)$$

It follows via the first item of (59) that the remainder as we have defined it does indeed have norm as small as possible. It follows via the second item of (59) that

$$N(\rho) \leq \frac{4}{7} N(\delta). \quad (64)$$

This represents a strengthening of the ordinary Euclidean condition (26) for this particular ring.

## 6. Reduced $\tau$ -adic NAF's

Having developed the modular reduction operation in  $\mathbb{Z}[\tau]$ , we now define the reduced  $\tau$ -adic NAF and apply it to the problem of efficient elliptic scalar multiplication.

### 6.1. Equivalence of $\tau$ -adic NAF's

We recall from the discussion at the end of §4.3 that our goal is a reduced  $\tau$ -adic NAF for  $n$ , equivalent to the ordinary  $\tau$ -adic NAF for  $n$  but only half as long.

We first define what we mean by *equivalent*. Let  $\mathcal{G}$  be a set of points on a Koblitz curve, and suppose that  $\gamma$  and  $\rho$  are two elements of  $\mathbb{Z}[\tau]$  for which  $\gamma P = \rho P$  for all  $P \in \mathcal{G}$ . Then we say that  $\text{TNAF}(\gamma)$  and  $\text{TNAF}(\rho)$  are *equivalent* with respect to  $\mathcal{G}$ . The terminology comes from the fact either  $\text{TNAF}(\gamma)$  or  $\text{TNAF}(\rho)$  can be used to multiply a point in  $\mathcal{G}$  by  $\gamma$ .

The following result of [18] gives a simple condition for two  $\tau$ -adic NAF's to be equivalent with respect to the entire set  $\mathcal{G} := E_a(\mathbb{F}_{2^m})$  of  $\mathbb{F}_{2^m}$ -rational points on  $E$ .

**PROPOSITION 65** *If  $\gamma$  and  $\rho$  are elements of  $\mathbb{Z}[\tau]$  with*

$$\gamma \equiv \rho \pmod{\tau^m - 1}, \quad (66)$$

*then*

$$\gamma P = \rho P$$

*for all  $P \in E_a(\mathbb{F}_{2^m})$ . Thus  $\text{TNAF}(\gamma)$  and  $\text{TNAF}(\rho)$  are equivalent with respect to  $E_a(\mathbb{F}_{2^m})$ .*

*Proof.* Applying the mapping (15)  $m$  times to  $P := (x, y)$ , we obtain

$$\tau^m P = (x^{2^m}, y^{2^m}).$$

But  $x^{2^m} = x$  and  $y^{2^m} = y$  since  $x$  and  $y$  are elements of  $\mathbb{F}_{2^m}$ . Thus

$$\tau^m P = P$$

for all  $P \in E_a(\mathbb{F}_{2^m})$ . It follows that

$$(\tau^m - 1)P = \mathcal{O} \quad (67)$$

for all  $P \in E_a(\mathbb{F}_{2^m})$ .

Suppose now that (66) holds. Then

$$\gamma = \rho + \kappa \cdot (\tau^m - 1)$$

for some  $\kappa \in \mathbb{Z}[\tau]$ . Therefore

$$\begin{aligned} \gamma P &= \rho P + \kappa \cdot (\tau^m - 1)P \\ &= \rho P + \kappa \mathcal{O} \\ &= \rho P + \mathcal{O} \\ &= \rho P, \end{aligned}$$

proving the result. ■

## 6.2. Equivalence for Points of Prime Order

One can sharpen (65) in the case of cryptographic interest, namely the main subgroup in a Koblitz curve of very nearly prime order (see §4.1). As in that section, the order of the curve is

$$\#E_a(\mathbb{F}_{2^m}) = f \cdot r,$$

where  $r$  is prime and  $f = 2$  or  $4$  according to (11). By (25), the element

$$\delta := (\tau^m - 1)/(\tau - 1) \quad (68)$$

has norm  $r$ .

In this section, we will show that one can weaken the hypotheses of (65) and still retain equivalence with respect to the main subgroup.

**PROPOSITION 69** *Let  $P$  be a point in the main subgroup in a Koblitz curve of very nearly prime order, and define  $\delta$  as in (68). Then*

$$\delta P = \mathcal{O}.$$

*Proof.* By (12), there is a point  $Q$  such that  $P = fQ$ . By (67), we have

$$(\tau^m - 1)Q = \mathcal{O}.$$

By (68), it follows that

$$\delta \cdot (\tau - 1)Q = \mathcal{O}. \quad (70)$$

Applying the operation  $\overline{\tau - 1}$  to both sides of (70), we obtain

$$\begin{aligned}\mathcal{O} &= \delta \cdot (\tau - 1) \cdot \overline{(\tau - 1)} Q \\ &= \delta \cdot N(\tau - 1) Q \\ &= \delta \cdot f Q \\ &= \delta P,\end{aligned}$$

proving the result. ■

**THEOREM 3** *Let  $P$  be a point in the main subgroup in an Koblitz curve of very nearly prime order, and define  $\delta$  as in (68). If  $\gamma$  and  $\rho$  are elements of  $\mathbb{Z}[\tau]$  with*

$$\gamma \equiv \rho \pmod{\delta},$$

then

$$\gamma P = \rho P.$$

Thus  $\text{TNAF}(\gamma)$  and  $\text{TNAF}(\rho)$  are equivalent with respect to the main subgroup.

*Proof.* The result follows from (69) in the same way as (65) follows from (67). ■

### 6.3. The Reduced $\tau$ -adic NAF

Suppose that  $E_a(\mathbb{F}_{2^m})$  has very nearly prime order, and that  $r$  is the order of the main subgroup. Let  $n$  be a positive integer less than  $r/2$ , and let  $\delta$  be as in (68). We define the *reduced  $\tau$ -adic NAF* of  $n$  to be

$$\text{RTNAF}(n) := \text{TNAF}(\rho),$$

where

$$\rho := n \bmod \delta.$$

It follows from Thm. 3 that  $\text{RTNAF}(n)$  and  $\text{TNAF}(n)$  are equivalent with respect to the main subgroup. Thus  $\text{RTNAF}(n)$  can be used in place of  $\text{TNAF}(n)$  for elliptic scalar multiplication in the main subgroup. It follows from the next theorem that this is a more efficient choice.

**THEOREM 4** *The average Hamming weight among reduced TNAF's is  $\sim m/3$ .*

*Proof.* Since the Hamming weight of an RTNAF is the product of its length and its density, we estimate both. We begin with the length.

It follows from (21) that

$$r = 2^{m-2+a} - (V_m - 1)2^{a-2}; \tag{71}$$

thus

$$r = 2^{m-2+a} + O(2^{m/2})$$

by (18). Now

$$N(\rho) \leq \frac{4}{7} r$$

by (64), so that

$$N(\rho) \leq \frac{2^{m+a}}{7} + O(2^{m/2}).$$

By (53), we may conclude that the length  $\ell_{\text{RTNAF}}$  of  $\text{RTNAF}(n)$  satisfies

$$\ell_{\text{RTNAF}} < m + a + .7082392.$$

Since  $\ell_{\text{RTNAF}}$  is an integer, it follows that

$$\ell_{\text{RTNAF}} \leq m + a. \tag{72}$$

We now consider the density of a RTNAF. B. Poonen has outlined a proof<sup>8</sup> that the TNAF's of integers modulo  $\tau^m - 1$  have average density  $1/3 + o(1)$  as  $m$  increases. The proof is easily modified to the case of RTNAF's. The result now follows via (72). ■

It follows from (54) that the RTNAF has about half the weight of the ordinary  $\tau$ -adic NAF. By (8), the weight of  $\text{RTNAF}(n)$  is about equal to that of  $\text{NAF}(n)$ . Thus, replacing  $\text{NAF}(n)$  by  $\text{RTNAF}(n)$  eliminates the elliptic doubles and keeps roughly constant the number of elliptic additions. We have therefore solved the difficulty mentioned at the end of §4.3.

## 7. Elliptic Scalar Multiplication on Koblitz Curves

We have now identified the procedure to use for elliptic scalar multiplication on a Koblitz curve, namely the analogue of the binary method using the RTNAF. We now present the explicit algorithms for computing the RTNAF, and give the elliptic scalar multiplication algorithms. Finally, we develop the  $\tau$ -adic analogue of the window method.

### 7.1. Computing the Reduced $\tau$ -adic NAF

To give an algorithm for computing the reduced  $\tau$ -adic NAF, we need only specialize the modular reduction algorithm (62) to the case of reducing an integer  $n$  modulo  $\delta$ .

Thus we set  $\gamma := n$ , and let  $\delta = d_0 + d_1 \tau$  be the norm- $r$  element given by (68). Then the integers  $g_i$  appearing in (62) are

$$g_i := s_i n,$$

where

$$s_0 := d_0 + \mu d_1$$

$$s_1 := -d_1.$$

The integers  $s_i$  can be expressed in terms of the Lucas sequence  $U_k$  via

$$s_i = \frac{(-1)^i}{f} (1 - \mu U_{m+3-a-i}), \quad (73)$$

where  $f$  is as in (11) and  $\mu$  is as in (14). Since the Lucas sequence  $U_k$  can be computed efficiently, so can the integers  $s_i$ . They need only be computed once per curve. Once that is done, the reduction method is as follows.

ROUTINE 74 (REDUCTION MODULO  $(\tau^m - 1)/(\tau - 1)$ )

*Per-Curve Parameters:*

$$m, a, s_0, s_1, r$$

*Input:*

$$n$$

*Output:*

$$\text{integers } r_0, r_1 \text{ specifying } r_0 + r_1 \tau := n \bmod (\tau^m - 1)/(\tau - 1)$$

*Computation:*

$$d_0 \leftarrow s_0 + \mu s_1$$

$$\lambda_0 \leftarrow s_0 n / r$$

$$\lambda_1 \leftarrow s_1 n / r$$

$$(q_0, q_1) \leftarrow \text{Round}(\lambda_0, \lambda_1) \quad (\text{via (60)})$$

$$r_0 \leftarrow n - d_0 q_0 - 2 s_1 q_1$$

$$r_1 \leftarrow s_1 q_0 - s_0 q_1$$

$$\text{Output } r_0, r_1$$

Note that one could store  $d_0$  rather than computing it during each reduction. This would save one integer addition per reduction, but require an additional  $\sim m/2$  bits of memory.

It will be helpful to have a geometric description of the the elements  $\rho$  resulting from reducing modulo  $\delta$  the integers  $n$  with  $0 \leq n < r$ . Following [4], we define the *Voronoi region*

$$\mathcal{V} := \{\delta \lambda : \lambda \in \mathcal{U}\}.$$

More explicitly,  $\mathcal{V}$  is given by the equations

$$-r \leq (2d_0 + \mu d_1)\lambda_0 + (\mu d_0 + 4d_1)\lambda_1 < r$$

$$-2r \leq (d_0 + 4\mu d_1)\lambda_0 - (3\mu d_0 - 2d_1)\lambda_1 < 2r$$

$$-2r \leq (d_0 - 3\mu d_1)\lambda_0 + (4\mu d_0 + 2d_1)\lambda_1 < 2r.$$

PROPOSITION 75 *The lattice points*

$$\lambda_0 + \lambda_1 \tau, \quad \lambda_i \in \mathbb{Z}$$

in  $\mathcal{V}$  are precisely the elements

$$n \bmod \delta, \quad 0 \leq n < r. \quad (76)$$

*Proof.* It follows from (63) that every element (76) must be a lattice point in  $\mathcal{V}$ . We must now prove the converse, *i.e.* that every lattice point in  $\mathcal{V}$  is the result of reducing modulo  $\delta$  an integer  $n$ .

We know that every lattice point in  $\mathcal{V}$  is the result of reducing modulo  $\delta$  an element of  $\mathbb{Z}[\tau]$ . Since the lattice points are incongruent modulo  $\delta$ , this means that they correspond to the elements of  $\mathbb{Z}[\tau]/\delta\mathbb{Z}[\tau]$ . Now, since  $\delta$  has norm  $r$ , then

$$\frac{\mathbb{Z}[\tau]}{\delta\mathbb{Z}[\tau]} \cong \frac{\mathbb{Z}}{r\mathbb{Z}}.$$

Thus there are precisely  $r$  lattice points in  $\mathcal{V}$ . Finally, since

$$\mathbb{Z} \cap \delta\mathbb{Z}[\tau] = r\mathbb{Z}$$

(see [23]), the integers  $0, 1, \dots, r-1$  are incongruent modulo  $\delta$ ; thus the  $r$  elements (76) are distinct. Thus each of the  $r$  lattice points in  $\mathcal{V}$  is an element (76) for some  $n$ . ■

Finally, we combine the modular reduction routine (74) with Alg. 1 to obtain the following algorithm for computing the reduced  $\tau$ -adic NAF.

ALGORITHM 2 (*Reduced  $\tau$ -adic NAF*)

*Parameters:*

$$m, a, s_0, s_1, r$$

*Input:*

*a positive integer  $n$*

*Output:*

RTNAF( $n$ )

*Computation:*

Compute  $\rho \leftarrow n \bmod \delta$  (via (74))  
 Compute  $\mathcal{S} \leftarrow \text{TNAF}(\rho)$  (via Alg. 1)  
 Output  $\mathcal{S}$

### 7.2. The $\tau$ -adic Method for Elliptic Scalar Multiplication

We now apply the reduced  $\tau$ -adic NAF to produce an efficient procedure for performing elliptic scalar multiplication on the main subgroup of  $E_a(\mathbb{F}_{2^m})$ . This is done by modifying Alg. 2 in the same way as was done for (4). This produces the following analogue of (6).

#### ALGORITHM 3 *Scalar Multiplication on Koblitz Curves*

*Per-Curve Parameters:*

$m, a, s_0, s_1, r$

*Input:*

$n$ , a positive integer less than  $r/2$

$P$ , a point in the main subgroup

*Output:*

$nP$

*Computation:*

```

Compute  $(r_0, r_1) \leftarrow n \bmod \delta$     (via (74))
Set  $Q \leftarrow \mathcal{O}$ 
 $P_0 \leftarrow P$ 
While  $r_0 \neq 0$  or  $r_1 \neq 0$ 
  If  $r_0$  odd then
    set  $u \leftarrow 2 - (r_0 - 2r_1 \bmod 4)$ 
    set  $r_0 \leftarrow r_0 - u$ 
    if  $u = 1$  then set  $Q \leftarrow Q + P_0$ 
    if  $u = -1$  then set  $Q \leftarrow Q - P_0$ 
  Set  $P_0 \leftarrow \tau P_0$     (=RightShift[ $P_0$ ])
  Set  $(r_0, r_1) \leftarrow (r_1 + \mu r_0/2, -r_0/2)$ 
EndWhile
Output  $Q$ 

```

Since  $\ell \approx m$ , then Alg. 3 requires  $\sim m/3$  additions and no doubles. This is at least 50% faster than any of the earlier versions, as shown in Table 1. The “length” and “density” columns in Table 1 give the approximate length of the relevant representation of the number and the average density of nonzero terms. The density figure of  $3/8$  for Koblitz’ “balanced” expansions is from experimental observation and may be only an approximation.

It should also be noted that the advantage of the three  $\tau$ -adic methods over the general methods is slightly overstated in Table 1, since doubling is more efficient than adding on a general elliptic curve.

Table 1. Comparison of elliptic scalar multiplication techniques.

Type of Curve	Method		Length of Expansion	Avg. Density	Avg. # of Elliptic Operations
General	Binary Method		$m$	$1/2$	$3m/2$
' '	Addition-Subtraction	(1989)	$m$	$1/3$	$4m/3$
Koblitz	Koblitz, Balanced	(1991)	$2m$	$3/8$	$3m/4$
' '	Meier-Staffelbach	(1992)	$m$	$1/2$	$m/2$
' '	$\tau$ -adic NAF	(1997)	$m$	$1/3$	$m/3$

### 7.3. The Width- $w$ $\tau$ -adic NAF

It remains to extend Alg. 3 to a “ $\tau$ -adic window method” analogous to (10). In order to do this, we first develop a width- $w$   $\tau$ -adic NAF.

In order to motivate our method, we recall the strategy used in the ordinary width- $w$  window method. With that method, we precompute and store the points  $uP$  as  $u$  runs over representatives of each odd congruence class (mod  $2^w$ ). When an odd integer  $c$  is encountered, the  $w$  rightmost bits are examined to determine which congruence class (mod  $2^w$ ) contains  $c$ . The corresponding point  $eP$  is subtracted off, and the new coefficient  $c - e$  is divisible by  $2^w$ .

For a  $\tau$ -adic window method, we precompute and store a point corresponding to each odd congruence class (mod  $\tau^w$ ). (By “odd” is meant that the element  $e_0 + e_1 \tau$  has odd  $e_0$ .) When an odd element  $r_0 + r_1 \tau$  is encountered, we must determine which congruence class (mod  $\tau^w$ ) contains  $r_0 + r_1 \tau$ . We can then subtract off the representative of that congruence class and produce a new coefficient which is divisible by  $\tau^w$ .

In analogy with the ordinary case, we can make the determination by examining the  $w$  rightmost bits of the appropriate combination of  $r_0$  and  $r_1$ . Define the quantities  $t_k$  via

$$t_k := 2 U_{k-1} U_k^{-1} \pmod{2^k}.$$

Since the Lucas elements  $U_k$  are odd, it follows that  $t_k$  is a well-defined integer modulo  $2^k$  that is even but not divisible by 4. Therefore

$$t_k^2 - \mu t_k + 2 \equiv 0 \pmod{2^k}$$

by (20). Thus  $t_k$  satisfies the same polynomial equation over  $\mathbb{Z}/2^k\mathbb{Z}$  that  $\tau$  satisfies over the complex numbers. It follows that the correspondence  $\tau \mapsto t_k$  induces a ring homomorphism from  $\mathbb{Z}[\tau]$  onto  $\mathbb{Z}/2^k\mathbb{Z}$  via

$$\begin{aligned} \phi_k: \mathbb{Z}[\tau] &\longrightarrow \mathbb{Z}/2^k\mathbb{Z} \\ u_0 + u_1 \tau &\longmapsto u_0 + u_1 t_k. \end{aligned}$$

It is easy to see that the odd elements of  $\mathbb{Z}[\tau]$  correspond precisely to the odd elements of  $\mathbb{Z}/2^k\mathbb{Z}$ .

We now compute the kernel of the mapping  $\phi_k$ .

LEMMA 77 *Let  $\alpha \in \mathbb{Z}[\tau]$ . Then  $\phi_k(\alpha) = 0$  if and only if  $\alpha$  is divisible by  $\tau^k$ .*

*Proof.* Since  $\phi_k(\tau) = t_k \equiv 2 \pmod{4}$ , it follows that  $\phi_k(\tau^j) = t_k^j$  is divisible by  $2^j$  but not by  $2^{j+1}$ . Thus the power of 2 dividing  $\phi_k(\alpha)$  is the power of  $\tau$  dividing  $\alpha$ . ■

It follows that each congruence class (mod  $\tau^k$ ) in  $\mathbb{Z}[\tau]$  corresponds under  $\phi_k$  with an element of  $\mathbb{Z}/2^k\mathbb{Z}$ . Moreover, the odd congruence classes (mod  $\tau^k$ ) in  $\mathbb{Z}[\tau]$  correspond under  $\phi_k$  with the odd elements of  $\mathbb{Z}/2^k\mathbb{Z}$ .

We now apply these results to construct a width- $w$  TNAF. It follows from the above correspondence that the odd numbers

$$\pm 1, \pm 3, \dots, \pm(2^{w-1} - 1)$$

are incongruent modulo  $\tau^w$ . Therefore, if

$$\alpha_u := u \bmod \tau^w,$$

then the numbers

$$\pm\alpha_1, \pm\alpha_3, \dots, \pm\alpha_{2^{w-1}-1}$$

also incongruent modulo  $\tau^w$ . If one precomputes the numbers  $\alpha_u$ , then there is the following simple algorithm for generating a width- $w$   $\tau$ -adic NAF.

ALGORITHM 4 (*Width- $w$   $\tau$ -adic NAF*)

*Parameters:*

$$a, w, t_w \\ \alpha_u := \beta_u + \gamma_u \tau \text{ for } u = 1, 3, \dots, (2^{w-1} - 1)$$

*Input:*

$$\text{an element } \rho := r_0 + r_1 \tau \text{ of } \mathbb{Z}[\tau]$$

*Output:*

$$\text{TNAF}_w(\rho)$$

*Computation:*

```

Set  $S \leftarrow \langle \rangle$ 
While  $r_0 \neq 0$  or  $r_1 \neq 0$ 
  If  $r_0$  odd
    then
      set  $u \leftarrow r_0 + r_1 t_w \bmod 2^w$ 
      if  $u > 0$ 
        then
           $\xi \leftarrow 1$ 
        else
           $\xi \leftarrow -1$ 

```

```

                 $u \leftarrow -u$ 
                set  $r_0 \leftarrow r_0 - \xi \beta_u$ 
                set  $r_1 \leftarrow r_1 - \xi \gamma_u$ 
                prepend  $\xi \alpha_u$  to  $\mathcal{S}$ 
            else
                prepend 0 to  $\mathcal{S}$ 
        Set  $(r_0, r_1) \leftarrow (r_1 + \mu r_0/2, -r_0/2)$ 
    EndWhile
    Output  $\mathcal{S}$ 

```

As before, we define the *reduced width- $w$   $\tau$ -adic NAF* of  $n$  to be

$$\text{RTNAF}_w(n) := \text{TNAF}_w(\rho),$$

where

$$\rho := n \bmod \delta.$$

The following algorithm uses the reduced width- $w$   $\tau$ -adic NAF to perform elliptic scalar multiplication.

ALGORITHM 5  *$\tau$ -adic Width- $w$  Window Method*

*Per-Curve Parameters:*

$$m, a, s_0, s_1, r, w, t_w$$

$$\alpha_u := \beta_u + \gamma_u \tau \text{ for } u = 1, 3, \dots, (2^{w-1} - 1)$$

*Input:*

*a positive integer  $n$*   
*an elliptic curve point  $P$*

*Output:*

*the point  $nP$*

*Precomputation:*

$$\text{Compute } P_u := \alpha_u P \text{ for } u = 1, 3, \dots, (2^{w-1} - 1)$$

*Computation:*

```

Set  $j \leftarrow 0$ 
Compute  $(r_0, r_1) := n \bmod \delta$  (via (74))
Set  $Q \leftarrow \mathcal{O}$ 
While  $r_0 \neq 0$  or  $r_1 \neq 0$ 
    If  $r_0$  odd
        then
            set  $u \leftarrow r_0 + r_1 t_w \bmod 2^w$ 
            if  $u > 0$ 

```

```

                                then
                                     $\xi \leftarrow 1$ 
                                else
                                     $\xi \leftarrow -1$ 
                                     $u \leftarrow -u$ 
                                set  $Q \leftarrow Q + \xi P_u$ 
Set  $j \leftarrow j + 1$ 
Set  $Q \leftarrow \tau^{-1}Q$            =LeftShift[ $Q$ ]
Set  $(r_0, r_1) \leftarrow (r_1 + \mu r_0/2, -r_0/2)$ 
EndWhile
Set  $Q \leftarrow \tau^j Q$            =RightShift[ $Q$ ] ( $j$  times)
Output  $Q$ 

```

Alg. 5 is a right-to-left window algorithm. We saw in §3.2 that this is usually impossible, but it can be done in the case of Koblitz curves because the Frobenius map  $\tau$  can be inverted efficiently.

#### 7.4. Performance

We now discuss the performance of Alg. 5.

The precomputation involves computing and storing  $2^{w-2} - 1$  points. For practical values of  $w$ , this requires  $2^{w-2} - 1$  elliptic additions and no storage beyond that needed to store the precomputed points. We illustrate for  $w = 5, a = 1$ .

In this case, the precomputation is to produce  $\alpha_u P$  for  $u = 3, 5, 7, \dots, 15$ . The numbers  $\alpha_u$  have the following TNAF's.

$$\begin{aligned} \alpha_3 &= \tau^2 - 1 \\ \alpha_5 &= \tau^2 + 1 \\ \alpha_7 &= -\tau^3 - 1 \\ \alpha_9 &= -\tau^5 - \tau^3 + 1 = -\tau^3 \alpha_5 + 1 \\ \alpha_{11} &= -\tau^4 - \tau^2 - 1 = -\tau^2 \alpha_5 - 1 \\ \alpha_{13} &= -\tau^4 - \tau^2 + 1 = -\tau^2 \alpha_5 + 1 \\ \alpha_{15} &= \tau^4 - 1 \end{aligned}$$

Therefore, the points  $P_u := \alpha_u P$  can be computed as follows.

$$\begin{aligned} P_3 &= \tau^2 P - P \\ P_5 &= \tau^2 P + P \\ P_7 &= -\tau^3 P - P \\ P_9 &= -\tau^3 P_5 + P \end{aligned}$$

Table 2. Performance at various widths.

Width	Number of Elliptic Operations		
	Precomputation	Real Time (avg)	Total (avg)
2	0	54.3	54.3
3	1	40.8	41.8
4	3	32.6	35.6
5	7	27.2	34.2
6	15	23.3	38.3
7	31	20.4	51.4

$$P_{11} = -\tau^2 P_5 - P$$

$$P_{13} = -\tau^2 P_5 + P$$

$$P_{15} = \tau^4 P - P$$

Since multiplication by  $\tau$  is essentially free, computing each point  $P_u$  requires one elliptic addition.

This approach works for all  $w$  from 3 to 8, and most likely for larger  $w$  as well.

The main computation is performed on a sequence of length  $\sim m$  and average density  $(w+1)^{-1}$ . It follows that a scalar multiplication on  $E_a(\mathbb{F}_{2^m})$  requires

$$\sim 2^{w-2} - 1 + \frac{m}{w+1} \quad (78)$$

elliptic additions on average.

Table 2 gives the performance of Alg. 5 on the curve  $E_1(\mathbb{F}_{2^{163}})$  for various widths, based on the estimate (78). (Entries are rounded to the nearest tenth of an integer.) The case  $w = 2$  is just Alg. 3. By choosing  $w = 5$ , one saves well over one-third the work. For larger  $w$ , the precomputation costs overshadow any savings on the real-time computation. Thus such larger widths would only be used for a long-term fixed point  $P$  (e.g. a public key).

It is remarkable that one can perform a general elliptic scalar multiplication on  $E_1(\mathbb{F}_{2^{163}})$  using only about 34 multiplicative inversions and 68 field multiplications.

One could obtain still further speedups by using more general window methods. These would be straightforward adaptations of existing methods such as those found in [13]. On the other hand, such methods are less automatic than the above fixed-width-window technique, so that more complicated up-front calculations are needed.

## 8. Efficient Modular Reduction

The modular reduction algorithm (74) is central to computing the reduced  $\tau$ -adic NAF, but it can be expensive to implement. The main computational difficulty is in the divisions

necessary in computing the expression (61). The purpose of this section is to present an equivalent method which is far more efficient.

### 8.1. Partial Modular Reduction

The equivalent method involves replacing the modular reduction process (74) by a simplified technique called *partial modular reduction*. As a result, one obtains an element  $\rho' \in \mathbb{Z}[\tau]$  that is congruent to  $n$  modulo  $\delta$ , but not necessarily of minimal norm.

This is done by replacing the rational numbers  $\lambda_i$  by approximations  $\lambda'_i$  which can be computed more efficiently. We denote by  $C$  the number of bits of accuracy of the approximations  $\lambda'_i$ . The larger  $C$  is, the more work is required in computing  $\lambda'_i$ , but the greater the likelihood that  $\lambda'_i$  will actually equal  $\lambda_i$ .

We begin by presenting the algorithm for computing  $\lambda'_i$ . We call this process *approximate division* because it replaces the division used to compute  $s_i n/r$  by two  $K$ -bit multiplications, where

$$K := \frac{m+5}{2} + C. \quad (79)$$

ALGORITHM 6 (*Approximate Division by  $r$* )

*Per-Curve Parameters:*

$$s_i, r \\ V_m \text{ (see §4.1)}$$

*Input:*

a positive integer  $n$  less than  $r/2$

*Output:*

$\lambda_i := s_i n/r$  to  $C$  bits of accuracy

*Computation:*

- 1  $n' \leftarrow \left\lfloor \frac{n}{2^{m-K-2+a}} \right\rfloor$
- 2  $g' \leftarrow s_i n'$
- 3  $h' \leftarrow \left\lfloor \frac{g'}{2^m} \right\rfloor$
- 4  $j' \leftarrow V_m h'$
- 5  $\ell' \leftarrow \text{Round} \left( \frac{g' + j'}{2^{K-C}} \right)$
- 6 Output  $\lambda' := \ell' / 2^C$ .

When the output of Alg. 6 is used in the modular reduction process (74), the resulting algorithm is what we call *partial reduction* modulo  $\delta := (\tau^m - 1)/(\tau - 1)$ . We shall write

$$\rho' = n \text{ partmod } \delta$$

for the result. The explicit algorithm is as follows.

ALGORITHM 7 (*Partial Reduction modulo  $(\tau^m - 1)/(\tau - 1)$* )

*Per-Curve Parameters:*

$$m, a, s_0, s_1, r$$

*Input:*

$$n$$

*Output:*

$$\text{integers } r'_0, r'_1 \text{ specifying } r'_0 + r'_1 \tau := n \text{ partmod } \delta$$

*Computation:*

$$d_0 \leftarrow s_0 + \mu s_1$$

$$\lambda'_0 \leftarrow s_0 n / r \text{ to } \sim K \text{ places via Alg. 6}$$

$$\lambda'_1 \leftarrow s_1 n / r \text{ to } \sim K \text{ places via Alg. 6}$$

$$(q'_0, q'_1) \leftarrow \text{Round}(\lambda'_0, \lambda'_1)$$

$$r'_0 \leftarrow n - d_0 q'_0 - 2 s_1 q'_1$$

$$r'_1 \leftarrow s_1 q'_0 - s_0 q'_1$$

$$\text{Output } r'_0, r'_1$$

We define the *partially reduced  $\tau$ -adic NAF* of  $n$  by

$$\text{PRTNAF}(n) := \text{TNAF}(\rho').$$

The element  $\rho'$  is congruent to  $n$  modulo  $\delta$ , since

$$\rho' = n - \kappa' \delta$$

where  $\kappa' := q'_0 + q'_1 \tau$ . Therefore,  $\text{PRTNAF}(n)$  and  $\text{TNAF}(n)$  are equivalent with respect to the main subgroup. Thus  $\text{PRTNAF}(n)$  can be used in place of  $\text{TNAF}(n)$  for elliptic scalar multiplication in the main subgroup.

On the other hand,  $\rho'$  may not be of minimal norm, as  $\rho$  would have been. Thus  $\text{PRTNAF}(n)$  may be longer than  $\text{RTNAF}(n)$ . If so, then more elliptic operations will be required in the scalar multiplication step. This would mitigate, and in all likelihood wipe out, the savings in avoiding the integer divisions in (74).

We examine this question in the next section. In particular, we shall prove that  $\text{PRTNAF}(n)$  is never much longer than  $\text{RTNAF}(n)$ . More importantly, we will show that the probability that  $\text{RTNAF}(n) \neq \text{PRTNAF}(n)$  can be made arbitrarily small, so that the possibility need not trouble us in practice.

## 8.2. Analysis of Partial Modular Reduction

We begin by analyzing the computational requirements of the partial modular reduction process.

**PROPOSITION 80** *Partial modular reduction, as implemented by Algs. 6 and 7, requires eight  $K$ -bit multiplications and no integer divisions, where  $K$  is given by (79).*

*Proof.* Since Alg. 7 requires two executions of Alg. 6 (one for each  $i$ ), a total of eight multiplications is required. No integer divisions are required since all denominators are powers of 2, so that the required divisions are implemented by simple bit shifts.

The only remaining issue is the number of bits in the numbers to be multiplied. Those numbers are  $n'$ ,  $h'$ ,  $V_m$ ,  $d_0$ , and  $s_i$  and  $q'_i$  for  $i = 0, 1$ . Each must be less than  $2^K$  in absolute value.

The desired bound on  $n'$  follows from the fact that  $n < r/2$ . The bounds for  $s_i$ ,  $d_0$ , and  $V_m$  come from (18) via (73). From these bounds is derived the bound for  $h'$ , along with

$$|\lambda'_i| < 2^{(m-1)/2}$$

for each  $i$ . Thus  $|q'_i| < 2^K$ . ■

We next address the accuracy of partial modular reduction.

**THEOREM 5** *If  $\lambda_i := g_i/r$ , and  $\lambda'_i$  is the approximation obtained by Alg. 6, then*

$$|\lambda_i - \lambda'_i| < 2^{-C}.$$

*Proof.* It follows from (71) that

$$\frac{g_i}{r} = \frac{g_i}{2^{m-a+2}} + \frac{(V_m - 1)g_i}{2^{2m-a+2}} + \frac{(V_m - 1)^2 g_i}{2^{2m} r},$$

so that

$$\frac{g_i}{r} = \frac{g_i}{2^{m-a+2}} + \frac{V_m g_i}{2^{2m-a+2}} + O(2^{-m/2}).$$

Therefore

$$\frac{2^K g_i}{r} = \frac{g_i}{2^{m-K-a+2}} \left( 1 + \frac{V_m}{2^m} \right) + O(2^{K-m/2}). \quad (81)$$

Now it follows from the definition of  $g'$  that

$$\left| \frac{g_i}{2^{m-K-a+2}} - g' \right| < |s_i|.$$

Thus by (81),

$$\frac{2^K g_i}{r} = g' + \frac{g' V_m}{2^m} + \theta_0 s_i + O(2^{K-m/2})$$

for some  $\theta_0$  with  $|\theta_0| < 1$ . It follows from the definition of  $j'$  that

$$\left| j' - \frac{g' V_m}{2^m} \right| < |V_m|.$$

Thus

$$\frac{2^K g_i}{r} = g' + j' + \theta_0 s_i + \theta_1 V_m + O(2^{K-m/2})$$

for some  $\theta_1$  with  $|\theta_1| < 1$ . It follows that

$$\frac{2^K g_i}{r} = g' + j' + \theta (|s_i| + |V_m|)$$

for some  $\theta$  with  $|\theta| < 1$ . Now

$$\begin{aligned} |s_i| + |V_m| &< 2^{1+m/2} + 2^{(m-1)/2} \\ &< 2^{(m+3)/2} \\ &\leq 2^{K-C-1}, \end{aligned}$$

so that

$$\left| \frac{2^K g_i}{r} - (g' + j') \right| < 2^{K-C-1}.$$

It follows that

$$\left| \frac{2^C g_i}{r} - \frac{g' + j'}{2^{K-C}} \right| < \frac{1}{2},$$

so that

$$\left| \frac{2^C g_i}{r} - \ell' \right| < 1.$$

Dividing this equation by  $2^C$  gives the result. ■

### 8.3. Worst-Case Effect of Partial Modular Reduction

In this section and the next, we retain the notation  $\kappa, \rho$  for the results of modular reduction and  $\kappa', \rho'$  for the results of partial modular reduction of  $n$  modulo  $\delta$ .

We have seen that

$$N(\rho) \leq \frac{4}{7} r,$$

and that therefore the TNAF( $\rho$ ) has length at most  $m + a$ . (See Thm. 4.) We now obtain corresponding results for  $\rho'$ .

LEMMA 82 *If  $\rho' := n \text{ partmod } \delta$ , then*

$$N(\rho') \leq \left( \frac{2}{\sqrt{7}} + \sqrt{N(\varepsilon)} \right)^2 r,$$

where  $\varepsilon = \kappa - \kappa'$ .

*Proof.* Since

$$\begin{aligned} \rho' - \rho &= (n - \kappa' \delta) - (n - \kappa \delta) \\ &= \varepsilon \delta, \end{aligned}$$

then

$$\sqrt{N(\rho')} \leq \sqrt{N(\rho)} + \sqrt{N(\varepsilon)} \sqrt{r}$$

by the Triangle Inequality (23). The result then follows from (64).  $\blacksquare$

THEOREM 6 *If  $C \geq 2$ , then the length of PRTNAF( $n$ ) is at most  $m + a + 3$ .*

*Proof.* By Thm. 5,

$$|\lambda_i - \lambda'_i| < \frac{1}{4}$$

for each  $i$ . It is easy to deduce that  $\lambda$  and  $\lambda'$  are either in the same copy of  $\mathcal{U}$  or adjacent copies. It follows that the centers of those copies satisfy

$$\kappa - \kappa' = 0, \pm 1, \pm \tau, \text{ or } \pm \bar{\tau}.$$

Thus  $N(\kappa - \kappa') \leq 2$ , so that

$$N(\rho') < 4.7095185 r$$

by (82). It follows via (53) that the length of TNAF( $\rho'$ ) is less than  $m + a + B$ , where  $B = 3.7511737$ . The result now follows from the fact that the length must be an integer.  $\blacksquare$

Thm. 6 states that the upper bound for the length of the PRTNAF is three bits higher than the corresponding bound for the RTNAF. In any given instance, the difference in length need not be three bits, since neither TNAF need attain its maximum possible length. Indeed, we will next prove that the two TNAF's are identical almost all of the time.

#### 8.4. Practical Effect of Partial Modular Reduction

PROPOSITION 83 *Let  $T > 0$  and  $\lambda := g/r$ , where  $g$  is an integer not divisible by  $r$  for which*

$$|g \bmod r| \leq \frac{r}{2} - T. \tag{84}$$

(The “mods” notation is as in §3.2.) If  $z$  is a rational number with

$$|z - \lambda| < \frac{T}{r},$$

then

$$\text{Round}(\lambda) = \text{Round}(z).$$

*Proof.* The inequality (84) can be rewritten as

$$|((\lambda))| \leq \frac{1}{2} - \frac{T}{r}.$$

If  $\kappa = \text{Round}(\lambda)$ , then,

$$|\lambda - \kappa| \leq \frac{1}{2} - \frac{T}{r}.$$

Therefore

$$\begin{aligned} |z - \kappa| &\leq |z - \lambda| + |\lambda - \kappa| \\ &< \frac{T}{r} + \left(\frac{1}{2} - \frac{T}{r}\right) \\ &\leq \frac{1}{2}, \end{aligned}$$

so that  $\text{Round}(z) = \kappa$ . ■

**COROLLARY 85** Let  $T > 0$ , and let  $s$  be an integer not divisible by  $r$ . If  $\varepsilon$  is a real number with

$$|\varepsilon| < \frac{T}{r},$$

then

$$\text{Round}\left(\frac{sn}{r}\right) = \text{Round}\left(\frac{sn}{r} + \varepsilon\right)$$

for all  $n$  between 0 and  $r/2$ , with at most  $T$  exceptions.

*Proof.* By (83), the only exceptions are the values of  $n$  for which

$$\frac{r}{2} - T \leq |sn \bmod r| < \frac{r}{2}. \quad \blacksquare$$

The following result is proved similarly to (85).

PROPOSITION 86 *Let  $T > 0$ , and let  $s$  be an integer not divisible by  $r$ . If*

$$|\varepsilon| < \frac{T}{r},$$

*then*

$$\left\lfloor \frac{sn}{r} \right\rfloor = \left\lfloor \frac{sn}{r} + \varepsilon \right\rfloor$$

*for all  $n$  between 0 and  $r/2$ , with at most  $T$  exceptions.*

LEMMA 87 *Let  $r$  be an odd prime, and let  $s_0$  and  $s_1$  be integers not divisible by  $r$ . For  $n$  an integer between 0 and  $r/2$ , let  $\lambda_i := s_i n/r$  for  $i = 0, 1$ . For each  $i$ , write*

$$\lambda_i = \lambda'_i + \delta_i$$

*where*

$$|\delta_i| < \frac{L}{r}. \tag{88}$$

*Then*

$$\text{Round}(\lambda_0, \lambda_1) = \text{Round}(\lambda'_0, \lambda'_1)$$

*for all values of  $n$  with at most  $14L$  exceptions.*

*Proof.* We examine the quantities  $f_i$ ,  $\eta_i$ , and  $h_i$  which are computed in the course of executing the algorithm (60) on  $\lambda_0 + \lambda_1 \tau$ , and compare them to the corresponding quantities  $f'_i$ ,  $\eta'_i$ , and  $h'_i$  from the same computation on  $\lambda'_0 + \lambda'_1 \tau$ .

It follows from (85) with  $T := L$  and  $\varepsilon := \delta_0$  that  $f_0 = f'_0$  for all  $n$  with at most  $L$  exceptions. Similarly,  $f_1 = f'_1$  for all  $n$  with at most  $L$  exceptions. Therefore, we have  $f_i = f'_i$  for both  $i$ , for all  $n$  with at most  $2L$  exceptions.

For each nonexceptional  $n$ , we have

$$\eta_i = \eta'_i + \delta_i$$

for each  $i$ . It follows that

$$\begin{aligned} |\eta - \eta'| &= |2\delta_0 + \mu \delta_1| \\ &< \frac{3L}{r}. \end{aligned}$$

It follows from (86) with  $T := 3L$  and  $\varepsilon := \eta - \eta'$  that

$$\lfloor \eta \rfloor = \lfloor \eta' \rfloor$$

for all  $n$  with at most  $3L$  exceptions. For these  $n$ ,

$$\eta \geq 1 \text{ if and only if } \eta' \geq 1$$

and

$$\eta < -1 \text{ if and only if } \eta' < -1.$$

Similarly, one finds that

$$\eta_0 - 3\mu\eta_1 \geq 1 \text{ if and only if } \eta'_0 - 3\mu\eta'_1 \geq 1,$$

$$\eta_0 - 3\mu\eta_1 < -1 \text{ if and only if } \eta'_0 - 3\mu\eta'_1 < -1$$

for all  $n$  with at most  $4L$  exceptions; and that

$$\eta_0 + 4\mu\eta_1 \geq 2 \text{ if and only if } \eta'_0 + 4\mu\eta'_1 \geq 2,$$

$$\eta_0 + 4\mu\eta_1 < -2 \text{ if and only if } \eta'_0 + 4\mu\eta'_1 < -2$$

for all  $n$  with at most  $5L$  exceptions.

It follows that all of the above conditions hold for all  $n$ , with at most  $14L$  exceptions. For each nonexceptional  $n$ , we have  $f_i = f'_i$  for both  $i$ , and all of the `if` conditions in (60) evaluate the same for the two computations. Thus  $h_i = h'_i$  for  $i = 0, 1$ . It follows that  $q_0 = q'_0$  and  $q_1 = q'_1$ , for each nonexceptional  $n$ .

**THEOREM 7** *If partial modular reduction with an accuracy of  $C$  bits is used, then the probability that  $\rho \neq \rho'$  is less than*

$$\text{Prob} < 2^{-(C-5)}.$$

*Proof.* Let  $L := 2^{-C}r$ . Then (88) holds by Thm. 5. By (87), the number of  $n$  for which  $\rho \neq \rho'$  is at most  $14L$ ; thus the probability that this happens is at most

$$\frac{14L}{r/2} < \frac{1}{2^{C-5}}. \quad \blacksquare$$

### 8.5. Conclusion

We now summarize the advantages and disadvantages of partial modular reduction.

If ordinary modular reduction (74) is used, then four  $K$ -bit multiplications are required, where  $K := (m + 3)/2$ . Also, two divisions are required, where the dividend is  $\sim 3m/2$  bits in length and the divisor is  $\sim m$  bits.

Using partial modular reduction (Algs. 6 and 7) with accuracy  $C$  has the following effects.

- The two divisions are no longer required.
- Four additional multiplications are required.
- The bit length of the integers to be multiplied is increased by  $C + 1$ .
- An answer  $\rho$  of non-minimal norm can be obtained with probability  $< 2^{-(C-5)}$ .

One can often choose  $C$  large enough that no non-minimal answer will be encountered in practice, at little practical cost. For example, if  $m = 163$ , then a multiplier of at least 83 bits is required. The likeliest available multiplier sizes are 96 or 128 bits. In the former case, one can choose  $C := 12$ , for a non-minimal answer at most once every 128 times. In that case, about three bits will be added to the TNAF, for an increase of about 1.8%. Thus the average cost comes to .08%, well worth it to avoid the divisions. For the latter case, one can choose  $C$  up to 54, guaranteeing that non-minimal answers will not occur in practice.

## 9. Recent Developments and Open Problems

Since the publication in the Proceedings of CRYPTO '97 of many of these results, the following advances have been made in the study of Koblitz curves.

- The results of [18] have recently been generalized [19] to curves defined over fields of  $2^d$  elements for small  $d$ . For example, the curves with complex multiplication by  $(\pm 1 + \sqrt{-15})/2$  are defined over  $\mathbb{F}_{2^2}$ . (The results of this paper should also carry over to this more general situation.)
- Analogues of Koblitz curves over fields of small odd characteristic have been studied [12], and representations analogous to the  $\tau$ -adic NAF (but with even fewer nonzero terms) have been obtained.
- It has been observed [3], [24] that the best square-root attack on elliptic curves [20] can be modified in the case of Koblitz curves. Rather than working with the points on the curve, one instead works with the cycles under the Frobenius operation. The resulting algorithm requires fewer steps than in the general case; however, each step is slightly more expensive.
- Finally, C. Günther and A. Stein have extended the concepts in this paper to the case of hyperelliptic curves [5].

The following are open problems in this field.

- We have presented worst-case upper bounds for the length of  $\text{TNAF}(\alpha)$  in terms of  $N(\alpha)$ . It may be more useful to possess analogous average-case results. For example, let

$$F(\alpha) := \ell(\alpha) - \log_2(N(\alpha))$$

for  $\alpha \in \mathbb{Z}[\tau]$ , where  $\ell(\alpha)$  denotes the length of  $\text{TNAF}(\alpha)$ . It follows from (53) that

$$-.5462682713 < F(\alpha) < 3.51559412$$

whenever  $\ell(\alpha) > 30$ . It would be a nice result to prove the existence of an average value for  $F(\alpha)$  and to evaluate it. (Note that the examples at the end of §4.3 prove that there is no asymptotic value for  $F(\alpha)$  as  $N(\alpha)$  gets large.)

- It is often required (*e.g.* in many cryptographic algorithms such as the ECDSA digital signature [8]) to take a random multiple of a point on an elliptic curve. More precisely, let  $P$  be in the main subgroup of a Koblitz curve of very nearly prime order. To take a random multiple of  $P$ , one generates a random integer  $n$  modulo  $r$ , reduces it modulo the element  $\delta$  given in (68), and uses the result  $\rho$  to compute  $\rho P$ . The techniques in this paper provide a particularly efficient way of doing the latter, but it would be cheaper still to produce a random point in the main subgroup directly. This could be done<sup>9</sup> by generating a “random  $\tau$ -adic NAF” of length  $m + a$ .

By a “random  $\tau$ -adic NAF” we mean a sequence of 0’s and  $\pm 1$ ’s which is generated as follows. Generate the first signed bit according to the following probability distribution:

$$u := \begin{cases} 0 & \text{prob} = 1/2 \\ 1 & \text{prob} = 1/4 \\ -1 & \text{prob} = 1/4. \end{cases} \quad (89)$$

To generate subsequent signed bits, follow each  $\pm 1$  by a 0, and generate each signed bit following a 0 according to (89).

The sequences generated in this way represent random selections from the set of all  $\tau$ -adic NAF’s of given length. In particular, each signed bit occurs with the proper average frequency.<sup>10</sup>

Once the sequence is generated, one computes  $\alpha P$  where  $\alpha$  is the element represented by the sequence.

This method gives random points, but their distribution is not known. It would be useful to measure how uniformly distributed such points are in the main subgroup.

More precisely, it follows from the proof of Thm. 4 that every  $n < r$  has a reduced  $\tau$ -adic NAF of length at most  $m + a$ . Thus every point in the main subgroup can be obtained by scalar multiplication using some nonadjacent sequence of  $m + a$  signed bits. It is easy to see that the number of such sequences is the integer closest to  $2^{m+a+2}/3$ . Since there are  $r \approx 2^{m-2+a}$  points in the main subgroup, the average number of sequences leading to a given point is  $16/3$ . It would be useful to know how much deviation there is from this average.

## A. Main Subgroup Membership

In this appendix we state and prove the conditions that are used in §4.1 to determine when a point on a Koblitz curve of very nearly prime order is an element of the main subgroup. We work on a general Weierstrass curve

$$E: y^2 + xy = x^3 + ax^2 + b$$

over  $GF(2^m)$ . It should be recalled that the order of  $E(GF(2^m))$  is always even, and divisible by 4 if and only if  $a$  has trace 0.

PROPOSITION 90 *Let  $(x_2, y_2)$  be a point on  $E$ . Then*

$$(x_2, y_2) = 2(x_1, y_1) \quad (91)$$

*for some  $(x_1, y_1)$  on  $E$  if and only if  $\text{Tr}(x_2) = \text{Tr}(a)$ .*

See [21].

PROPOSITION 92 *Suppose that  $a$  has trace 0. Let  $(x_4, y_4)$  be a point on  $E$ . Then*

$$(x_4, y_4) = 4(x_1, y_1) \quad (93)$$

*for some  $(x_1, y_1)$  on  $E$  if and only if  $\text{Tr}(x_4) = 0$  and*

$$\text{Tr}(y_4) = \text{Tr}(\lambda x_4) \quad (94)$$

*for some  $\lambda$  satisfying*

$$\lambda^2 + \lambda = x_4 + a. \quad (95)$$

*Proof.* Suppose first that (93) holds for some  $(x_1, y_1)$ . Let

$$(x_2, y_2) := 2(x_1, y_1),$$

so that  $(x_4, y_4) = 2(x_2, y_2)$ . By (90), it follows that

$$\text{Tr}(x_4) = \text{Tr}(x_2) = 0. \quad (96)$$

Let

$$\lambda := x_2 + \frac{y_2}{x_2}; \quad (97)$$

then it follows from the doubling formula that (95) holds and that

$$y_4 := x_2^2 + (\lambda + 1)x_4. \quad (98)$$

Since  $\text{Tr}(x_2^2) = \text{Tr}(x_2)$ , it follows (via (96)) that (94) holds.

Conversely, suppose now that  $\text{Tr}(x_4) = 0$ , and that (94) holds for some  $\lambda$  for which (95) holds. Then in fact (94) holds for *either* value of  $\lambda$  satisfying (95). It follows from (90) that  $(x_4, y_4) = 2(x_2, y_2)$  for some  $(x_2, y_2)$  on  $E$ . We may conclude that (97) holds (where we have replaced  $\lambda$  by  $\lambda + 1$  if necessary). Moreover, (98) holds, so that

$$y_4 + \lambda x_4 = x_2^2 + x_4.$$

It follows that  $\text{Tr}(x_2^2) = 0$ , so that  $\text{Tr}(x_2) = 0$ . By (90), it follows that  $(x_2, y_2) = 2(x_1, y_1)$  for some  $(x_1, y_1)$  on  $E$ . Thus  $(x_4, y_4) = 4(x_1, y_1)$ . ■

## B. Gaussian Normal Bases

In anticipation of the specific Koblitz curves presented in Appendix 9, we describe the Gaussian normal basis representation for the binary field  $GF(2^m)$ .

Let  $m$  be prime, and let  $T$  be an even positive integer for which  $p := Tm + 1$  is prime. Then a Type  $T$  Gaussian normal basis exists for  $GF(2^m)$  if and only if  $m$  is relatively prime to  $(p - 1)/k$ , where  $k$  is the order of 2 modulo  $p$ . For each prime  $m$ , a Type  $T$  Gaussian normal basis exists for some  $T$ . To maximize efficiency, the lowest available value of  $T$  is used.

Once the type  $T$  has been identified, the multiplication rule can be constructed. One first constructs a function  $F(\underline{u}, \underline{v})$  on inputs

$$\underline{u} = (u_0 u_1 \dots u_{m-1})$$

and

$$\underline{v} = (v_0 v_1 \dots v_{m-1})$$

as follows. Let  $u$  be an integer having order  $T$  modulo  $p$ . Compute the sequence  $J(1), J(2), \dots, J(p - 1)$  as follows:

```

1           Set  $w \leftarrow 1$ 
2           For  $j$  from 0 to  $T - 1$  do
2.1         Set  $n \leftarrow w$ 
2.2         For  $i$  from 0 to  $m - 1$  do
2.2.1       Set  $J(n) \leftarrow i$ 
2.2.2       Set  $n \leftarrow 2n \bmod p$ 
2.3         Set  $w \leftarrow uw \bmod p$ 

```

Then  $F$  is given by the formula

$$F(\underline{u}, \underline{v}) := \sum_{k=1}^{p-2} u_{J(k+1)} v_{J(p-k)}.$$

This computation need only be performed once per basis. Given the function  $F$  for  $\mathcal{B}$ , one computes the product

$$(c_0 c_1 \dots c_{m-1}) = (a_0 a_1 \dots a_{m-1}) \times (b_0 b_1 \dots b_{m-1})$$

as follows.

```

1           Set  $(u_0 u_1 \dots u_{m-1}) \leftarrow (a_0 a_1 \dots a_{m-1})$ 
2           Set  $(v_0 v_1 \dots v_{m-1}) \leftarrow (b_0 b_1 \dots b_{m-1})$ 
3           For  $k$  from 0 to  $m - 1$  do
3.1         Compute  $c_k := F(\underline{u}, \underline{v})$ 
3.2         Set  $\underline{u} \leftarrow \text{LeftShift}(\underline{u})$  and  $\underline{v} \leftarrow \text{LeftShift}(\underline{v})$ ,
              where LeftShift denotes the
              circular left shift operation.
4           Output  $c := (c_0 c_1 \dots c_{m-1})$ 

```

*Example:* For the type 4 normal basis for  $GF(2^7)$ , one has  $p = 29$  and  $u = 12$  or  $17$ . Thus the values of  $F$  are given by

$$\begin{array}{llll} F(1) = 0 & F(8) = 3 & F(15) = 6 & F(22) = 5 \\ F(2) = 1 & F(9) = 3 & F(16) = 4 & F(23) = 6 \\ F(3) = 5 & F(10) = 2 & F(17) = 0 & F(24) = 1 \\ F(4) = 2 & F(11) = 4 & F(18) = 4 & F(25) = 2 \\ F(5) = 1 & F(12) = 0 & F(19) = 2 & F(26) = 5 \\ F(6) = 6 & F(13) = 4 & F(20) = 3 & F(27) = 1 \\ F(7) = 5 & F(14) = 6 & F(21) = 3 & F(28) = 0 \end{array}$$

Therefore

$$\begin{aligned} F(\underline{u}, \underline{v}) &= u_0 v_1 + u_1 (v_0 + v_2 + v_5 + v_6) + u_2 (v_1 + v_3 + v_4 + v_5) \\ &\quad + u_3 (v_2 + v_5) + u_4 (v_2 + v_6) + u_5 (v_1 + v_2 + v_3 + v_6) \\ &\quad + u_6 (v_1 + v_4 + v_5 + v_6). \end{aligned}$$

Thus, if

$$a = (1010111) \quad \text{and} \quad b = (1100001),$$

then

$$\begin{aligned} c_0 &= F((1010111), (1100001)) = 1, \\ c_1 &= F((0101111), (1000011)) = 0, \\ &\quad \vdots \\ c_6 &= F((1101011), (1110000)) = 1, \end{aligned}$$

so that  $c = ab = (1011001)$ . □

### C. Standard Curves

The following five Koblitz curves appear in the document "Recommended Elliptic Curves for Federal Government Use," issued July 1999 by NIST and available on their website

<http://csrc.nist.gov/encryption>.

Each curve has very nearly prime order. For each curve is given a base point  $G = (G_x, G_y)$  generating the main subgroup.

*Curve K-163*

$$a = 1$$

$$r = 5846006549323611672814741753598448348329118574063$$

$$\text{Polynomial Basis} \quad t^{163} + t^7 + t^6 + t^3 + 1$$



Curve K-409

$a = 0$

$r = 33052798439512429947595765401638551991420234148214060964232\backslash$   
 $4395022880711289249191050673258457777458014096366590617731\backslash$   
 $358671$

Polynomial Basis  $t^{409} + t^{87} + 1$

$G_x =$		060f05f	658f49c1	ad3ab189
	0f718421	0efd0987	e307c84c	27accfb8
	c460189e	b5aaaa62	ee222eb1	b35540cf
$G_y =$		1e36905	0b7c4e42	acba1dac
	bf04299c	3460782f	918ea427	e6325165
	da5f6c42	e9c55215	aa9ca27a	5863ec48
			d8e0286b	

Type 4 Normal Basis

$G_x =$		1b559c7	cba2422e	3affe133
	43e808b5	5e012d72	6ca0b7e6	a63aeafb
	10ca0fcf	98350c3b	7f89a975	4a8e1dc0
$G_y =$		16d8c42	052f07e7	713e7490
	eff318ba	1abd6fef	8a5433c8	94b24f5c
	852496fb	ee803a47	bc8a2038	78ebf1c4
			99afd7d6	

Curve K-571

$a = 0$

$r = 19322687615086291723476759454659936721494636648532174993286\backslash$   
 $1762572575957114478021226813397852270671183470671280082535\backslash$   
 $1461273674974066617311929682421617092503555733685276673$

Polynomial Basis  $t^{571} + t^{10} + t^5 + t^2 + 1$

$G_x =$		26eb7a8	59923fbc	82189631
	f8103fe4	ac9ca297	0012d5d4	60248048
	43709584	93b205e6	47da304d	b4ceb08c
	494776fb	988b4717	4dca88c7	e2945283
$G_y =$		349dc80	7f4fbf37	4f4aeade
	3bca9531	4dd58cec	9f307a54	ffc61efc
	9d4979c0	ac44aea7	4fbegbb9	f772aedc
	7ba7af1b	320430c8	591984f6	01cd4c14
				3ef1c7a3

Type 10 Normal Basis

$G_x =$		04bb2db	a418d0db	107adae0
	03427e5d	7cc139ac	b465e593	4f0bea2a
	c29b3d5b	9aa7a1fd	fd5d8be6	6057c100
				8e71e484

$$G_y = \begin{array}{cccccc} \text{bcd98f22} & \text{bf847642} & \text{37673674} & \text{29ef2ec5} & \text{bc3ebcf7} & \\ & & \text{44cbb57} & \text{de20788d} & \text{2c952d7b} & \\ \text{56cf39bd} & \text{3e89b189} & \text{84bd124e} & \text{751ceff4} & \text{369dd8da} & \\ \text{c6a59e6e} & \text{745df44d} & \text{8220ce22} & \text{aa2c852c} & \text{fcbbef49} & \\ \text{ebaa98bd} & \text{2483e331} & \text{80e04286} & \text{feaa2530} & \text{50caff60} & \end{array}$$

## Acknowledgments

The author wishes to thank Dan Gordon, Neal Koblitz, Julio Lopez, and Bob Reiter for many helpful comments and suggestions.

## Notes

1. We restrict our attention to elliptic curves that are not *supersingular*, since such curves are cryptographically weaker than ordinary curves [16]. But see [12] for cryptographic applications of supersingular curves.
2. This does not cause confusion, because the origin is never on  $E$ .
3. There do exist special-purpose improvements to the basic elliptic operations, *e.g.* [14], but they are not relevant to this paper.
4. It is easy to prove there is no left-to-right method for computing the NAF. On the other hand, there exist signed binary expansions that are as good as the NAF and that can be computed from left to right.
5. More elaborate window methods exist (see [4]), but they can require a great deal of initial calculation and seldom do much better than the technique presented here.
6. See, *e.g.* [6].
7. As R. Schroepel has remarked, these algorithms are also useful when using a polynomial basis, since squaring is still relatively efficient in that case.
8. A brief summary of Poonen's approach is given in [4] by D. Gordon, who has since presented a more detailed version of the proof.
9. This is an adaptation of an idea of H. Lenstra (see [10]).
10. It was proved in (32) that 0 occurs with frequency  $2/3$  after the initial  $\pm 1$ . It is easy to see that 1 and  $-1$  are equally likely on average.

## References

1. D. W. Ash, I. F. Blake, and S. Vanstone, Low complexity normal bases, *Discrete Applied Math.*, Vol. 25 (1989) pp. 191–210.
2. E. Berlekamp, *Algebraic Coding Theory*, Aegean Park Press (1984).
3. R. Gallant, R. Lambert, and S. Vanstone, Improving the parallelized Pollard lambda search on binary anomalous curves, *Math. Comp.*, posted on May 19, 1999, PII: S-0025-5718(99)01119-9 (to appear in print).
4. D. Gordon, A survey of fast exponentiation methods, *J. Algs.*, Vol. 27 (1998) pp. 129–146.
5. C. Günther and A. Stein (to appear).
6. Institute of Electrical and Electronics Engineers, *IEEE P1363: Standard Specifications for Public-Key Cryptography*, Draft 10 (1999).
7. T. Itoh, O. Teechai, and S. Trojii, A fast algorithm for computing multiplicative inverses in  $GF(2^t)$ , *J. Soc. Electron. Comm. (Japan)*, Vol. 44 (1986) pp. 31–36.
8. D. Johnson and A. Menezes, *The Elliptic Curve Digital Signature Algorithm (ECDSA)*, Univ. of Waterloo (1999) <http://cacr.math.uwaterloo.ca>

9. D. E. Knuth, *Seminumerical Algorithms*, Addison-Wesley (1981).
10. N. Koblitz, CM curves with good cryptographic properties, *Proc. Crypto '91*, Springer-Verlag (1992) pp. 279–287.
11. N. Koblitz, *A Course of Number Theory and Cryptography*, 2nd ed., Springer-Verlag (1994).
12. N. Koblitz, An elliptic curve implementation of the Finite Field Digital Signature Algorithm, *Proc. Crypto '98*, Springer-Verlag (1998) pp. 327–337.
13. K. Koyama and Y. Tsuruoka, Speeding up elliptic cryptosystems by using a signed binary window method, *Proc. Crypto '92*, Springer-Verlag (1993) pp. 345–357.
14. J. Lopez, Fast multiplication on elliptic curves over  $GF(2^m)$  without precomputation (preprint).
15. F. Morain and J. Olivos, Speeding up the computations on an elliptic curve using addition-subtraction chains, *Inform. Theor. Appl.*, Vol. 24 (1990) pp. 531–543.
16. A. Menezes, T. Okamoto and S. Vanstone, Reducing elliptic curve logarithms to logarithms in a finite field, *IEEE Transactions on Information Theory*, Vol. 39 (1993) pp. 1639–1646.
17. A. Menezes, P. van Oorschot, and S. Vanstone, *Handbook of Applied Cryptography*, CRC Press (1997).
18. W. Meier and O. Staffelbach, Efficient multiplication on certain non-supersingular elliptic curves, *Proc. Crypto '92*, Springer-Verlag (1993) pp. 333–344.
19. V. Müller, Fast multiplication on elliptic curves over small fields of characteristic two, *J. Crypt.*, Vol. 11 (1998) pp. 219–234.
20. P. van Oorschot and M. Weiner, Parallel collision search with cryptanalytic applications, *J. Crypt.*, Vol. 12 (1999) pp. 1–28.
21. G. Seroussi, Compact representations of elliptic curve points over  $GF(2^m)$ , <http://grouper.ieee.org/groups/1363/contributions/hp.ps>
22. J. Silverman, *The Arithmetic of Elliptic Curves*, Springer-Verlag (1986).
23. I. Stewart and D. Tall, *Algebraic Number Theory*, 2nd. ed., Chapman and Hall (1987).
24. M. Weiner and R. Zuccherato, Faster attacks on elliptic curve cryptosystems, *Selected Areas in Cryptography*, Springer-Verlag (1999) pp. 190–200.