

Efficient and Secure Algorithms for GLV-Based Scalar Multiplication and their Implementation on GLV-GLS Curves

Armando Faz-Hernández^{1*}, Patrick Longa², and Ana H. Sánchez³

¹ Institute of Computing,
University of Campinas, Brazil
`armfazh@ic.unicamp.br`

² Microsoft Research,
One Microsoft Way, Redmond, USA
`plonga@microsoft.com`

³ Computer Science Department, CINVESTAV-IPN, México
`asanchez@computacion.cs.cinvestav.mx`

Abstract. We propose efficient algorithms and formulas that improve the performance of *side-channel protected* scalar multiplication exploiting the Gallant-Lambert-Vanstone (CRYPTO 2001) and Galbraith-Lin-Scott (EUROCRYPT 2009) methods. Firstly, by adapting Feng et al.’s recoding to the GLV setting, we derive new regular algorithms for variable-base scalar multiplication that offer protection against simple side-channel and timing attacks. Secondly, we propose an efficient technique that interleaves ARM-based and NEON-based multiprecision operations over an extension field, as typically found on GLS curves and pairing computations, to improve performance on modern ARM processors. Finally, we showcase the efficiency of the proposed techniques by implementing a state-of-the-art GLV-GLS curve in twisted Edwards form defined over \mathbb{F}_{p^2} , which supports a four dimensional decomposition of the scalar and runs in constant time, i.e., it is fully protected against timing attacks. For instance, using a precomputed table of only 512 bytes, we compute a variable-base scalar multiplication in 92,000 cycles on an Intel Ivy Bridge processor and in 244,000 cycles on an ARM Cortex-A15 processor. Our benchmark results and the proposed techniques contribute to the improvement of the state-of-the-art performance of elliptic curve computations. Most notably, our techniques allow us to reduce the cost of adding protection against timing attacks in the GLV-based variable-base scalar multiplication computation to below 10%.

Keywords. Elliptic curves, scalar multiplication, side-channel protection, constant-time computation, GLV method, GLS method, GLV-GLS curve, x64 processor, ARM processor, NEON instructions.

1 Introduction

Let P be a point of prime order r on an elliptic curve over \mathbb{F}_p containing a degree-2 endomorphism ϕ . The Gallant-Lambert-Vanstone (GLV) method computes the scalar multiplication kP as $k_1P + k_2\phi(P)$ [15]. If k_1, k_2 have approximately half the bitlength of the original scalar k , one should expect an elimination of half the number of doublings by using the Straus-Shamir simultaneous multi-scalar multiplication technique. Thus, the method is especially useful for speeding up the case in which the base point P is variable, known as variable-base scalar multiplication. Later, Galbraith et al. [14] showed how to exploit the Frobenius endomorphism to enable the use of the GLV approach on a wider set of curves defined over the quadratic extension field \mathbb{F}_{p^2} . Since then, significant research has been performed to improve the performance [29, 23] and to explore the applicability to other settings [19, 34] or to higher dimensions on genus one curves [23, 30] and genus two curves [7, 8]. Unfortunately, most of the work and comparisons with other approaches have been carried out with *unprotected* algorithms and implementations. In fact, little effort has been done to investigate methods for protecting GLV-based implementations against side-channel attacks. Just recently, Longa and Sica [30] used the regular windowed recoding by Okeya and Takagi [33] in combination with interleaving [15, 32] to make their four-dimensional implementation constant time. However, the use of this

* Author became affiliated to University of Campinas at the time of publication.

standard approach in the GLV paradigm incurs in a high cost in terms of storage and computing performance because of the high number of required precomputations. This issue worsens for higher dimensions [8].

In this work, we propose a new signed representation, called GLV-based Sign-Aligned Column (GLV-SAC), that gives rise to a new method for scalar multiplication using the GLV method. We depart from the traditional approach based on interleaving or joint sparse form and adapt the recoding by Feng et al. [11], originally intended for standard comb-based fixed-base scalar multiplication, to the computation of GLV-based variable-base scalar multiplication. The method supports a regular execution and thus provides a first layer of protection against some simple side-channel (SSCA) attacks such as simple power analysis (SPA) [26]. Moreover, it does not require dummy operations, making it resilient to safe-error attacks [40, 41], and can be used as a basis for constant-time implementations secure against timing attacks [25, 9, 2, 35]. In comparison with the best previous approaches, the method improves the computing performance, especially during the potentially expensive precomputation stage, and allows us to save *at least* half of the storage requirement for precomputed values without impacting performance. For instance, the method injects a 17% speedup in the overall computation and a 78% reduction in the memory consumption for a GLV-GLS curve using a 4-GLV decomposition (see §5). The savings in memory without impacting performance are especially relevant for the deployment of GLV-based implementations in constrained devices. Depending on the cost of endomorphisms, the improvement provided by the method is expected to increase for higher-degree decompositions.

Processors based on the ARM architecture are widely used in modern smartphones and tablets due to their low power consumption. The ARM architecture comes equipped with 16 32-bit registers and an instruction set including 32-bit operations, which in most cases can be executed in one cycle. To boost performance in certain applications, some ARM processors include a powerful set of vector instructions known as NEON. This consists of a 128-bit Single Instruction Multiple Data (SIMD) engine that includes 16 128-bit registers. Recent research has exploited NEON to accelerate cryptographic operations [6, 18, 36]. On one hand, the interleaving of ARM and NEON instructions is a well-known technique (with increasing potential on modern processors) that can be exploited in cryptography; e.g., see [6]. On the other hand, the vectorized computation using NEON can be advantageously exploited to compute independent multiplications, as found in operations over \mathbb{F}_{p^2} ; e.g., see [36]. In this work, we take these optimizations further and propose a technique that interleaves ARM-based and NEON-based multiprecision operations, such as multiplication, squaring and modular reduction, in extension field operations in order to maximize the inherent parallelism and hide the execution latency. The technique is especially relevant for implementing the quadratic extension field layer, as found in GLS curves [14] and pairing computations [1]. For instance, it injects a significant speedup in the range 17%-34% in the scalar multiplication execution on the targeted GLV-GLS curve (see §4 and §5).

To demonstrate the efficiency of our techniques, we implement the state-of-the-art twisted Edwards GLV-GLS curve over \mathbb{F}_{p^2} with $p = 2^{127} - 5997$, recently proposed by Longa and Sica [30]. This curve, referred to as **Ted127-g1v4**, supports a 4-GLV decomposition. Moreover, we also present efficient algorithms for implementing field and quadratic extension field operations targeting our 127-bit prime on x64 and ARM platforms. We combine and exploit incomplete reduction [39] and lazy reduction [38], expanding techniques by [29]. These optimized operations are then applied to state-of-the-art twisted Edwards formulas [3, 22] to speed up computations in the setting of curves over \mathbb{F}_{p^2} . Our implementations of variable-base scalar multiplication target modern x64 and ARM processors, and include *full* protection against timing attacks: the scalar is decomposed and recoded (in constant time) in a regular pattern using the proposed GLV-SAC representation, secret-data conditional branches are avoided and memory accesses (over precomputed points) are performed in constant time.

Notably, we show that the proposed algorithms and formulas reduce dramatically the cost of protecting against timing attacks and the storage for precomputations, and set a new speed record for protected software. For instance, a protected variable-based elliptic curve scalar multiplication on curve **Ted127-g1v4** runs in 96,000 cycles on an Intel Sandy Bridge, using only 512 bytes of memory for precomputed values. This is 30% faster, using almost 1/5 of the storage, than a previous implementation by Longa and Sica [30] also based on curve **Ted127-g1v4** that computes the same operation in 137,000 cycles using 2.25KB of memory for precomputations. Moreover, this result is only 5% slower, using 1/2 of the storage, than the *unprotected* computation by the same authors, which runs in 91,000 cycles and uses 1KB of memory. This not only

represents a new speed record for protected software but also marks the first time that a constant-time variable-base scalar multiplication computation is performed under 100K cycles on an Intel processor. Similar results are obtained for ARM processors exploiting the technique that interleaves NEON and ARM-based operations.

This paper is organized as follows. In §2, we give some preliminaries about the GLV and GLS methods, and side-channel attacks. In §3, we present the new GLV-based representation and the corresponding scalar multiplication method. We describe the implementation of curve `Ted127-glv4` as well as optimized algorithms for field, extension field and point operations targeting x64 and ARM platforms in §4. In this section, we also discuss the interleaving technique for ARM. Finally, in §5, we perform an analysis of the proposed methods and present benchmark results of scalar multiplication on several x64 and ARM processors.

2 Preliminaries

2.1 The GLV and GLS Methods

In this section, we briefly describe the GLV and GLS methods in a generic, m dimensional framework. Let C be a curve defined over a finite field \mathbb{F}_p equipped with an efficiently computable endomorphism ϕ . The GLV method to compute scalar multiplication [15] consists of first decomposing the scalar k into sub-scalars k_i for $0 \leq i < m$ and then computing $\sum_{i=0}^{m-1} k_i D_i$ using the Straus-Shamir trick for simultaneous multi-scalar multiplication, where D_0 is the input divisor from the divisor class group of the curve and $D_i = \phi^i(D_0)$. If all of the sub-scalars have approximately the same bitlength, the number of required doublings is reduced to approximately $\log_2 r/m$, where r is the prime order of the curve subgroup. Special curves equipped with endomorphisms which are different to the Frobenius endomorphism are known as GLV curves.

The GLS method [14, 13] lifts the restriction to special curves and exploits an endomorphism ψ arising from the p -power Frobenius endomorphism on a wider set of curves C' defined over an extension field \mathbb{F}_{p^k} that are \mathbb{F}_{p^n} -isogenous to curves C/\mathbb{F}_p , where $k|n$. Equipped with ψ to perform the scalar decomposition, one then proceeds to apply the GLV method as above. More complex decompositions arise by applying the GLS paradigm to GLV curves (a.k.a. GLV-GLS curves [14, 30]).

These techniques have received lots of attention recently, given their significant impact in the performance of curve-based systems. Longa and Gebotys [29] report efficient implementations of GLS curves over \mathbb{F}_{p^2} using two-dimensional decompositions. In [23], Hu, Longa and Xu explore a GLV-GLS curve over \mathbb{F}_{p^2} supporting a four-dimensional decomposition. In [7], Bos et al. study two and four-dimensional decompositions on genus 2 curves over \mathbb{F}_p . Bos et al. [8] explore the combined GLV-GLS approach over genus 2 curves defined over \mathbb{F}_{p^2} , which supports an 8-GLV decomposition. In the case of binary GLS elliptic curves, Oliveira et al. [34] report the implementation of a curve exploiting the 2-GLV method. More recently, Guillevic and Ionica [17] show how to exploit the 4-GLV method on certain genus one curves defined over \mathbb{F}_{p^2} and genus two curves defined over \mathbb{F}_p ; and Smith [37] proposes a new family of elliptic curves that support 2-GLV decompositions.

From all these works, only [30] and [34] include side-channel protection in their GLV-based implementations.

2.2 Side-Channel Attacks and Countermeasures

Side-channel attacks [25] exploit leakage information obtained from the physical implementation of a cryptosystem to get access to private key material. Examples of physical information that can be exploited are power, time, electromagnetic emanations, among others. In particular, much attention has been put on timing [25, 9] and simple power attacks (SPA) [26], given their broad applicability and relatively low costs to be realized in practice. Traditionally, the different attacks can also be distinguished by the number of traces that are exploited in the analysis: simple side-channel attacks (SSCA) require only one trace (or very few traces) to observe the leakage that directly reveals the secret bits, whereas differential side-channel attacks (DSCA) require many traces to perform a statistical analysis on the data. The feasibility of these attacks depends on the targeted application, but it is clear that SSCA attacks are feasible in a wider range of scenarios. In this work, we focus on methods that minimize the risk posed by SSCA attacks such as SPA, and timing attacks.

In curve-based cryptosystems, the first step to achieve protection against these attacks is to use regular algorithms for performing scalar multiplication (other methods involve the use of unified formulas, but these are generally expensive). One efficient approach in this direction is to recode the scalar to a representation exhibiting a regular pattern. In particular, for the case of variable-base scalar multiplication, the regular windowed recoding proposed by Okeya and Takagi [33] and further analyzed by Joye and Tunstall [24] represents one of the most efficient alternatives. Nevertheless, in comparison with the standard width- w non-adjacent form (w NAF) [20] used in unprotected implementations, the Okeya-Takagi recoding increases the nonzero density from $1/(w+1)$ to $1/(w-1)$. In contrast, side-channel protected methods for scalar multiplication exploiting the GLV method have not been fully studied. Furthermore, we note that methods typically efficient in the standard case are not necessarily efficient in the GLV paradigm. For example, in [30], Longa and Sica apply the Okeya-Takagi recoding to protect scalar multiplication on a GLV-GLS curve using a four-dimensional GLV decomposition against timing attacks. The resulting protected implementation is about 30% more expensive than the unprotected version. In this work, we aim at reducing that gap, providing efficient methods that can be exploited to improve and protect GLV and GLS-based implementations.

The comb method [27] is an efficient approach for the case of fixed-base scalar multiplication. However, in its original form, the method is unprotected against SSCA and timing attacks. An efficient approach to achieve a regular execution is to recode the scalar using signed nonzero representations such as LSB-set [11], MSB-set [12] or SAB-set [21]. A key observation in this work is that the basic version of the fixed-base comb execution (i.e., without exploiting multiple tables) has several similarities with a GLV-based variable-base execution. So it is therefore natural to adapt these techniques to the GLV setting to achieve side-channel protection. In particular, the LSB-set representation is a good candidate, given that an analogue of this method in the GLV setting minimizes the cost of precomputation.

2.3 The Least Significant Bit - Set (LSB-Set) Representation

Feng et al. [11] proposed a clever signed representation, called LSB-set, that is based on the equivalence $1 \equiv 1\bar{1} \dots \bar{1}$ (assuming the notation $-1 \equiv \bar{1}$), and used it to protect the comb method [27] in the computation of fixed-base scalar multiplication (we refer to this method as LSB-set comb scalar multiplication). Next, we briefly describe the LSB-set recoding and its application to fixed-base scalar multiplication. The reader is referred to [27] and [11] for complete details about the original comb method and the LSB-set comb method, respectively.

Let t be the bitlength of a given scalar k . Assume that k is partitioned in w consecutive parts of $d = \lceil t/w \rceil$ bits each, padding k with $(dw - t)$ zeros to the left. Let the updated binary representation of k be $(k_{l-1}, k_{l-2}, \dots, k_0)$, where $l = dw$. One can visualize the bits of k in matrix form by considering the w pieces as the rows and arranging them from top to bottom. The LSB-set recoding consists of first applying the transformation $1 \mapsto 1\bar{1} \dots \bar{1}$ to the least significant d bits of the scalar (i.e., the first row in the matrix) and, then, converting every bit k_i in the remaining rows in such a way that output digits b_i for $d \leq i \leq (l-1)$ are in the digit set $\{0, b_{i \bmod d}\}$. That is, digits in the same column are 0 or share the same sign. Then, for computing a comb fixed-base scalar multiplication, one scans the “digit-columns” in the matrix from left to right. Since every digit-column is nonzero by definition, the execution consists of a doubling-addition computation at every iteration, which provides protection against certain SSCA attacks such as SPA.

3 The GLV-Based Sign-Aligned Column (GLV-SAC) Representation

In this section, we introduce a variant of the LSB-set recoding that is amenable for the computation of side-channel protected variable-base scalar multiplication in the GLV setting. The new recoding is called GLV-Based Sign-Aligned Column (GLV-SAC). Also, we present a new method for GLV-based scalar multiplication exploiting the proposed representation.

In the following, we first discuss the GLV-SAC representation in a generic setting. In Section 3.2, we discuss variants that are expected to be more efficient when $m = 2$ and $m \geq 8$. To simplify the descriptions,

we assume in the remainder that we are working on an elliptic curve. The techniques and algorithms can be easily extended to other settings such as genus 2 curves.

Let $\{k_0, k_1, \dots, k_j, \dots, k_{m-1}\}$ be a set of positive sub-scalars in the setting of GLV with dimension m . The basic idea of the new recoding is to have one of the sub-scalars of the m -GLV decomposition, say k_J , represented in signed nonzero form and acting as a “sign-aligner”. The latter means that k_J determines the sign of all the digits of remaining sub-scalars according to their relative position.

The GLV-SAC representation has the following properties:

- (i) The length of the digit representation of every sub-scalar k_j is fixed and given by $l = \lceil \log_2 r/m \rceil + 1$, where r is the prime subgroup order.
- (ii) Exactly one sub-scalar, which should be odd, is expressed by a signed nonzero representation $k_J = (b_{l-1}^J, \dots, b_0^J)$, where all digits $b_i^J \in \{1, -1\}$ for $0 \leq i < l$.
- (iii) All the sub-scalars k_j , with exception of k_J from (ii), are expressed by signed representations $(b_{l-1}^j, \dots, b_0^j)$ such that $b_i^j \in \{0, b_i^J\}$ for $0 \leq i < l$.

In the targeted setting, (i) and (ii) guarantee a constant-time execution regardless of the value of the scalar k and without having to appeal to masking for dealing with the identity element. Item (iii) allows us to reduce the size of the precomputed table by a factor of 2, while minimizing the cost of precomputation.

Note that we do not impose any restriction on which sub-scalar should be designated as k_J . In some settings, choosing any of the k_j (with the exception of the one corresponding to the base point, i.e., k_0) could lead to the same performance in the precomputation phase and be slightly faster than $k_J = k_0$, if one takes into consideration the use of mixed point additions. The condition that k_J should be odd enables the conversion of any integer to a full signed nonzero representation using the equivalence $1 \equiv 1\bar{1} \dots \bar{1}$. To deal with this restriction during the scalar multiplication, we first convert the selected sub-scalar k_J to odd (if even), and then make the corresponding correction at the end (more details can be found in Section 3.1). Finally, the reader should note that the GLV-SAC representation, in the way we describe it above, assumes that the sub-scalars are all positive. This restriction is imposed in order to achieve the minimum length $l = \lceil \log_2 r/m \rceil + 1$ in the representation. Note that it is possible to lift this restriction if needed in a certain setting (the analysis of this case is included in the extended paper version [10]).

An efficient algorithm to recode the sub-scalars to GLV-SAC proceeds as follows. Assume that each sub-scalar k_j is padded with zeros to the left until reaching the fixed length $l = \lceil \log_2 r/m \rceil + 1$, where r is the prime order of the curve subgroup. After choosing a suitable k_J to act as the “sign-aligner”, the sub-scalar k_J is recoded to signed nonzero digits b_i^J using the equivalence $1 \equiv 1\bar{1} \dots \bar{1}$. Remaining sub-scalars are then recoded in such a way that output digits at position i are in the set $\{0, b_i^J\}$, i.e., nonzero digits at the same relative position share the same sign. This is shown as Algorithm 1.

Algorithm 1 Protected Recoding Algorithm for the GLV-SAC Representation.

Input: m l -bit positive integers $k_j = (k_{l-1}^j, \dots, k_0^j)_2$ for $0 \leq j < m$, an odd “sign-aligner” $k_J \in \{k_j\}^m$, where $l = \lceil \log_2 r/m \rceil + 1$, m is the GLV dimension and r is the prime group order.

Output: $(b_{l-1}^j, \dots, b_0^j)_{\text{GLV-SAC}}$ for $0 \leq j < m$, where $b_i^j \in \{1, -1\}$, and $b_i^j \in \{0, b_i^J\}$ for $0 \leq j < m$ and $j \neq J$.

```

1:  $b_{l-1}^J = 1$ 
2: for  $i = 0$  to  $(l - 2)$  do
3:    $b_i^J = 2k_{i+1}^J - 1$ 
4: for  $j = 0$  to  $(m - 1), j \neq J$  do
5:   for  $i = 0$  to  $(l - 1)$  do
6:      $b_i^j = b_i^J \cdot k_0^j$ 
7:    $k_j = \lfloor k_j/2 \rfloor - \lfloor b_i^j/2 \rfloor$ 
8: return  $(b_{l-1}^j, \dots, b_0^j)_{\text{GLV-SAC}}$  for  $0 \leq j < m$ .
```

We highlight that, in contrast to [11, Alg. 4] and [12, Alg. 2], our recoding algorithm is simpler and exhibits a regular and constant time execution, making it resilient to timing attacks. Moreover, Algorithm 1

can be implemented very efficiently by exploiting the fact that the only purpose of the recoded digits from the sub-scalar k_J is, by definition, to determine the sign of their corresponding digit-columns (see details in Alg. 2 below). Since $k_{i+1}^J = 0$ and $k_{i+1}^J = 1$ indicate that the corresponding output digit-column i will be negative and positive, respectively, Step 3 of Algorithm 1 can be reduced to $b_i^J = k_{i+1}^J$ by assuming the convention $b_i^J = 0$ to indicate negative and $b_i^J = 1$ to indicate positive, for $0 \leq i < l$. Following this convention, further efficient simplifications are possible for Steps 6 and 7.

3.1 GLV-Based Scalar Multiplication using GLV-SAC

We now present a new method for computing variable-base scalar multiplication using the GLV method and the GLV-SAC representation (see Algorithm 2). To simplify the description, we assume that k_0 is fixed as the “sign-aligner” k_J (it is easy to modify the algorithm to set any other sub-scalar to k_J). The basic idea is to arrange the sub-scalars, after being converted to their GLV-SAC representation, in matrix form from top to bottom, with sub-scalar $k_J = k_0$ at the top, and then run a simultaneous multi-scalar multiplication execution scanning digit-columns from left to right. By using the GLV-SAC recoding, every digit-column i is expected to be nonzero and have any of the possible combinations $[b_i^{m-1}, \dots, b_i^2, b_i^1, b_i^0]$, where $b_i^0 \in \{1, -1\}$, and $b_i^j \in \{0, b_i^0\}$ for $1 \leq j < m$. Since nonzero digits in the same column have the same sign, one only needs to precompute all the positive combinations $P_0 + u_1P_1 + \dots + u_{m-1}P_{m-1}$ with $u_j \in \{0, 1\}$, where P_j are the base points of the sub-scalars. Assuming that negation of group elements is inexpensive in a given curve subgroup, negative values can be computed on-the-fly during the evaluation stage.

Algorithm 2 Protected m -GLV Variable-Base Scalar Multiplication using the GLV-SAC Representation.

Input: Base point P_0 of order r and $(m-1)$ points P_j for $1 \leq j < m$ corresponding to the endomorphisms, m scalars $k_j = (k_{t_j-1}^j, \dots, k_0^j)_2$ for $0 \leq j < m$, $l = \lceil \frac{\log_2 r}{m} \rceil + 1$ and $\max(t_j) = \lceil \frac{\log_2 r}{m} \rceil$.

Output: kP .

Precomputation stage:

1: Compute $P[u] = P_0 + u_0P_1 + \dots + u_{m-2}P_{m-1}$ for all $0 \leq u < 2^{m-1}$, where $u = (u_{m-2}, \dots, u_0)_2$.

Recoding stage:

2: $\text{even} = k_0 \bmod 2$

3: **if** $\text{even} = 0$ **then** $k_0 = k_0 - 1$

4: Pad each k_j with $(l-t_j)$ zeros to the left for $0 \leq j < m$ and convert them to the GLV-SAC representation using Algorithm 1 s.t. $k_j = (b_{l-1}^j, \dots, b_0^j)_{\text{GLV-SAC}}$. Set digit-columns $\mathbb{K}_i = [b_i^{m-1}, \dots, b_i^2, b_i^1] \equiv [b_i^{m-1}2^{m-2} + \dots + b_i^22 + b_i^1]$ and digit-column signs $s_i = b_i^0$ for $0 \leq i \leq l-1$.

Evaluation stage:

5: $Q = s_{l-1}P[\mathbb{K}_{l-1}]$

6: **for** $i = l-2$ **to** 0 **do**

7: $Q = 2Q$

8: $Q = Q + s_iP[\mathbb{K}_i]$

9: **if** $\text{even} = 0$ **then** $Q = Q + P_0$

10: **return** Q

Since the GLV-SAC recoding requires that the “sign-aligner” k_J (in this case, k_0) be odd, k_0 is subtracted by one if it is even in Step 3 of Algorithm 2. The correction is then performed at the end of the evaluation stage at Step 9. These computations, as well as the accesses to the precomputed table, should be performed in constant time to guarantee protection against timing attacks. For example, in the implementation discussed in Section 5, the value $P[\mathbb{K}_i]$ required at Step 8 is retrieved from memory by performing a linear pass over the whole precomputed table using conditional move instructions. The final value $s_iP[\mathbb{K}_i]$ is then obtained by performing a second linear pass over the points $P[\mathbb{K}_i]$ and $-P[\mathbb{K}_i]$. Similarly, to realize Step 9, we always carry out the computation $Q' = Q + P_0$ and then perform a linear pass over the points Q and Q' using conditional move instructions to transfer the correct value to the final destination.

Note that Algorithm 2 assumes a decomposed scalar as input. This is sufficient in some settings, in which randomly generated sub-scalars could be provided. However, in others settings, one requires

to calculate the sub-scalars in a decomposition phase. We remark that this computation should also be computed in constant time for protecting against timing attacks (e.g., see the details for `Ted127-glv4` in §5).

Example 1. Let $m = 4$, $\log_2 r = 16$ and $kP = 11P_0 + 6P_1 + 14P_2 + 3P_3$. Using Algorithm 1, the corresponding GLV-SAC representation with fixed length $l = \lceil 16/4 \rceil + 1 = 5$ is given by (arranged in matrix form from top to bottom as required in Alg. 2)

$$\begin{bmatrix} k_0 \\ k_1 \\ k_2 \\ k_3 \end{bmatrix} \equiv \begin{bmatrix} 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix} \equiv \begin{bmatrix} 1 & \bar{1} & 1 & \bar{1} & 1 \\ 1 & \bar{1} & 0 & \bar{1} & 0 \\ 1 & 0 & 0 & \bar{1} & 0 \\ 0 & 0 & 1 & \bar{1} & 1 \end{bmatrix}$$

According to Algorithm 2, digit columns are given by $\mathbb{K}_0 = [100] = 4$, $\mathbb{K}_1 = [\bar{1}\bar{1}\bar{1}] = 7$, $\mathbb{K}_2 = [100] = 4$, $\mathbb{K}_3 = [00\bar{1}] = 1$ and $\mathbb{K}_4 = [011] = 3$, and their corresponding s_i are $s_0 = 1$, $s_1 = -1$, $s_2 = 1$, $s_3 = -1$ and $s_4 = 1$. Precomputed values $P[u]$ are given by $P[0] = P_0$, $P[1] = P_0 + P_1$, $P[2] = P_0 + P_2$, $P[3] = P_0 + P_1 + P_2$, $P[4] = P_0 + P_3$, $P[5] = P_0 + P_1 + P_3$, $P[6] = P_0 + P_2 + P_3$ and $P[7] = P_0 + P_1 + P_2 + P_3$. At Step 5 of Alg. 2, we compute $Q = s_4 P[\mathbb{K}_4] = P[3] = P_0 + P_1 + P_2$. The main loop in the evaluation stage is then executed as follows

i	3	2	1	0
$2Q$	$2P_0 + 2P_1 + 2P_2$	$2P_0 + 2P_1 + 4P_2$	$6P_0 + 4P_1 + 8P_2 + 2P_3$	$10P_0 + 6P_1 + 14P_2 + 2P_3$
$Q + s_i P[\mathbb{K}_i]$	$P_0 + P_1 + 2P_2$	$3P_0 + 2P_1 + 4P_2 + P_3$	$5P_0 + 3P_1 + 7P_2 + P_3$	$11P_0 + 6P_1 + 14P_2 + 3P_3$

Cost Analysis. To simplify comparisons, we will only consider a setting in which precomputed points are left in some projective system. When converting points to affine is convenient, one should include the cost of this conversion. Also, we do not consider optimizations exploiting cheap endomorphism mappings during precomputation, since this is dependent on a specific application. The reader is referred to Section 5 for a more precise comparison in a practical implementation using a GLV-GLS twisted Edwards curve.

The cost of the proposed m -GLV variable-base scalar multiplication using the GLV-SAC representation (Alg. 2) is given by $(l - 1)$ doublings and l additions during the evaluation stage using 2^{m-1} points, where $l = \lceil \frac{\log_2 r}{m} \rceil + 1$. Naively, precomputation costs $2^{m-1} - 1$ additions (in practice, several of these additions might be performed using cheaper mixed additions). So the total cost is given by $(l - 1)$ doublings and $(l + 2^{m-1} - 1)$ additions.

In contrast, the method based on the regular windowed recoding [33] used in [30] requires $(l - 1)$ doublings and $m \cdot (l - 1)/(w - 1) + 2m - 1$ additions during the evaluation stage and m doublings with $m \cdot (2^{w-2} - 1)$ additions during the precomputation stage, using $m \cdot (2^{w-2} + 1)$ points (naive approach without exploiting endomorphisms). If, for example, $r = 256$, $m = 4$ and $w = 5$ (typical parameters to achieve 128-bit security on a curve similar to `Ted127-glv4`), the new method costs 64 doublings and 72 additions using 8 points, whereas the regular windowed method costs 68 doublings and 99 additions using 36 points. Thus, the new method improves performance while reduces dramatically the number of precomputations (in this case, to almost 1/5 of the storage). Assuming that one addition costs 1.3 doublings, the expected speedup is 20%.

Certainly, one can reduce the number of precomputations when using the regular windowed recoding by only precomputing multiples corresponding to one or some of the sub-scalars. However, these savings in memory come at the expense of computing endomorphisms during the evaluation stage, which can cost from several multiplications [7] to approximately one full point addition each (see Appendix A). The proposed method always requires the minimum storage without impacting performance.

The basic GLV-SAC representation and its corresponding scalar multiplication are particularly efficient for four-dimensional GLV. In the following section, we discuss variants that are efficient for $m = 2$ and $m \geq 8$.

3.2 Windowed and Partitioned GLV-SAC: Case of Dimension 2 and ≥ 8

In some cases, the performance of the proposed scalar multiplication can be improved further by combining windowed techniques with the GLV-SAC recoding. Given a window width w , assume that a set of sub-scalars k_j has been padded with enough zeros to the left to guarantee that $w|l$, where l is the expected length of an extended GLV-SAC representation that we refer to as w GLV-SAC. The basic idea is to join every w consecutive digits in the w GLV-SAC representation, and precompute all possible values $P[u] = u'P_0 + u_0P_1 + \dots + u_{m-2}P_{m-1}$ for $0 \leq u < 2^{wm-1}$ and $u' \in \{1, 3, \dots, 2^w - 1\}$ (again, points corresponding to negative values of u' can be computed on-the-fly). Scalar multiplication then proceeds by scanning w -digit columns from left to right.

Conveniently, Algorithm 1 can also be used to obtain w GLV-SAC(k_j), with the only change in the fixed length to $l = (\lceil \log_2 r/w \rceil + 1) + (\lceil \log_2 r/w \rceil + 1) \bmod w$.

Example 2. Let $m = 2, \log_2 r = 8, w = 2$ and $kP = 11P_0 + 14P_1$. Using Algorithm 1, the corresponding w GLV-SAC representation with fixed length $l = \lceil 8/2 \rceil + 1 + (\lceil 8/2 \rceil + 1) \bmod 2 = 6$, arranged in matrix form from top to bottom, is given by

$$\begin{bmatrix} k_0 \\ k_1 \end{bmatrix} \equiv \begin{bmatrix} 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 & 0 \end{bmatrix} \equiv \begin{bmatrix} 1 & \bar{1} & \bar{1} & 1 & \bar{1} & 1 \\ 1 & \bar{1} & 0 & 0 & \bar{1} & 0 \end{bmatrix} \quad (1)$$

The 2-digit columns are given by $\mathbb{K}_0 = [\bar{2}\bar{1}] = 3, \mathbb{K}_1 = [0\bar{1}] = 1$ and $\mathbb{K}_2 = [11] = 2$, and their corresponding s_i are $s_0 = -1, s_1 = -1$ and $s_2 = 1$. Precomputed values $P[u]$ are given by $P[0] = P_0 - P_1, P[1] = P_0, P[2] = P_0 + P_1, P[3] = P_0 + 2P_1, P[4] = 3P_0, P[5] = 3P_0 + P_1, P[6] = 3P_0 + 2P_1$ and $P[7] = 3P_0 + 3P_1$. In the evaluation stage we first compute $Q = s_2P[\mathbb{K}_2] = P[2] = P_0 + P_1$ and then execute

i	1	0
$2^w Q$	$4P_0 + 4P_1$	$12P_0 + 16P_1$
$Q + s_i P[\mathbb{K}_i]$	$3P_0 + 4P_1$	$11P_0 + 14P_1$

Since the requirement of precomputations, given by 2^{wm-1} , increases rapidly as w and m grow, windowed GLV-SAC is especially attractive for 2-GLV implementations. In this case, by fixing $w = 2$ the number of precomputed points is only 8. At the same performance level (in the evaluation stage), this is approximately half the memory requirement of a method based on the regular windowed recoding [33]¹.

Whereas joining columns in the representation matrix is amenable for small m using windowing, for large m it is recommended to join rows instead. We illustrate the approach with $m = 8$. Given a set of sub-scalars k_j for $0 \leq j < 8$, we first partition it in c consecutive sub-sets k'_i such that $c|8$, and then convert every sub-set to the GLV-SAC representation (using Alg. 1). In this case, every column in the matrix consists of c sub-columns, each one corresponding to a sub-set k'_i . Scalar multiplication then proceeds by scanning c sub-columns per iteration from right to left. Thus, with this “partitioned” GLV-SAC approach, one increases the number of point additions per iteration in the main loop of Alg. 2 from one to c . However, the number of required precomputations is reduced from 2^{m-1} to $c \cdot 2^{\frac{m}{c}-1}$. For example, for $m = 8$, this approach reduces the number of points from 128 to only 16 if c is fixed to 2 (each sub-table corresponding to a sub-set of scalars contains 8 points). At the same performance level (in the evaluation stage), this is approximately half the memory requirement of a method based on the regular windowed recoding [33], as discussed by the recent work by Bos et al. [8]. Performance is also expected to improve since the number of point operations for precomputation is significantly reduced. Note that, if one only considers positive sub-scalars and the endomorphism mapping is inexpensive in comparison to point addition, then sub-tables can be computed by simply applying the endomorphism to the first sub-table arising from the base point P_0 . In some instances, such as the 8-GLV in [8], this approach is expected to reduce further the cost of precomputation. Although an

¹ However, in some cases one can afford the reduction of precomputations from 16 to 8 when using the windowed recoding if endomorphisms are cheap and can be computed on-the-fly during the evaluation stage; e.g., see [34].

issue arises when sub-scalars can also be negative, this can be dealt with by adding one extra bit containing the sign to the representation. We give the full details in the extended paper version [10].

4 High-Speed Implementation on GLV-GLS Curves

In this section, we describe implementation aspects of the GLV-GLS curve **Ted127-glv4**. We present optimized algorithms for prime field, extension field and point arithmetic. We also present the technique of interleaving NEON and ARM-based multiprecision operations over \mathbb{F}_{p^2} . Although our techniques are especially tuned for the targeted curve, we remark that they can be adapted and exploited in other scenarios.

4.1 The Curve

For complete details about the four-dimensional method using GLV-GLS curves, the reader is referred to [14] and [31]. We use the following GLV-GLS curve in twisted Edwards form [30], referred to as **Ted127-glv4**:

$$E'_{TE}/\mathbb{F}_{p^2} : -x^2 + y^2 = 1 + dx^2y^2, \quad (2)$$

where \mathbb{F}_{p^2} is defined as $\mathbb{F}_p[i]/(i^2 - \beta)$, $\beta = -1$ is a quadratic non-residue in \mathbb{F}_p and $u = 1 + i$ is a quadratic non-residue in \mathbb{F}_{p^2} . Also, $p = 2^{127} - 5997$, $d = 170141183460469231731687303715884099728 + 116829086847165810221872975542241037773i$ and $\#E'_{TE}(\mathbb{F}_{p^2}) = 8r$, where r is the 251-bit prime $2^{251} - 255108063403607336678531921577909824432295$. E'_{TE} is isomorphic to the Weierstrass curve $E'_W/\mathbb{F}_{p^2} : y^2 = x^3 - 15/2 u^2x - 7u^3$, which is the quadratic twist of a curve isomorphic to the GLV curve $E_W/\mathbb{F}_p : y^2 = 4x^3 - 30x - 28$ (see [30, Section 5]). E'_{TE}/\mathbb{F}_{p^2} is equipped with two efficiently computable endomorphisms Φ and Ψ defined over \mathbb{F}_{p^2} , which enable a four-dimensional decomposition for any scalar $k \in [1, r-1]$ in the subgroup generated by a point P of order r and, consequently, enable a four-dimensional scalar multiplication given by

$$kP = k_1P + k_2\Phi(P) + k_3\Psi(P) + k_4\Psi\Phi(P), \quad \text{with } \max_i(|k_i|) < Cr^{1/4}$$

where $C = 179$ [30]. Let $\zeta_8 = u/\sqrt{2}$ be a primitive 8th root of unity. The affine formulas for $\Phi(x, y)$ and $\Psi(x, y)$ are given by

$$\Phi(x, y) = \left(-\frac{(\zeta_8^3 + 2\zeta_8^2 + \zeta_8)xy^2 + (\zeta_8^3 - 2\zeta_8^2 + \zeta_8)x}{2y}, \frac{(\zeta_8^2 - 1)y^2 + 2\zeta_8^3 - \zeta_8^2 + 1}{(2\zeta_8^3 + \zeta_8^2 - 1)y^2 - \zeta_8^2 + 1} \right) \quad \text{and} \quad \Psi(x, y) = \left(\zeta_8x^p, \frac{1}{y^p} \right),$$

respectively. It can be verified that $\Phi^2 + 2 = 0$ and $\Psi^2 + 1 = 0$. The formulas in homogeneous projective coordinates can be found in Appendix A.

Note that **Ted127-glv4** has $a = -1$ (in the twisted Edwards equation; see [3]), corresponding to the most efficient set of formulas proposed by Hisil et al. [22]. Although GLV-GLS curves with suitably chosen parameters when transformed to twisted Edwards form offer roughly the same performance, as discussed in [30], there are certain differences in the cost of formulas for computing the endomorphisms Φ and Ψ . Curve **Ted127-glv4** exhibits relatively efficient formulas for computing the endomorphisms in comparison with other GLV-GLS curves from [30]. On the other hand, our selection of the pseudo-Mersenne prime $p = 2^{127} - 5997$ enables efficient field arithmetic by exploiting lazy and incomplete reduction techniques (see the next section for details). Also, since $p \equiv 3 \pmod{4}$, -1 is a quadratic non-residue in \mathbb{F}_p , which minimizes the cost of multiplication over \mathbb{F}_{p^2} by transforming multiplications by β to inexpensive subtractions.

4.2 Field Arithmetic

For field inversion, we use the modular exponentiation $a^{p-2} \pmod{p} \equiv a^{-1}$ using a fixed and short addition chain. This method is simple to implement and is naturally protected against timing attacks.

In the case of a pseudo-Mersenne prime of the form $p = 2^m - c$, with c small, field multiplication can be efficiently performed by computing an integer multiplication followed by a modular reduction exploiting the

special form of the prime. This separation of operations also enables the use of lazy reduction in the extension field arithmetic. For x64, integer multiplication is implemented in product scanning form (a.k.a Comba’s method), mainly exploiting the powerful 64-bit unsigned multiplier instruction. Let $0 \leq a, b < 2^{m+1}$. To exploit the extra room of one bit in our targeted prime $2^{127} - 5997$, we first compute $M = a \cdot b = 2^{m+1}M_H + M_L$ followed by the reduction step $R = M_L + 2cM_H \leq 2^{m+1}(2c+1) - 2$. Then, given $R = 2^m R_H + R_L$, we compute $R_L + cR_H \pmod{p}$, where $R_L, cR_H < 2^m$. This final operation can be efficiently carried out by employing the modular addition proposed by Bos et al. [7] to get the final result in the range $[0, p - 1]$. Note that the computation of field multiplication above naturally accepts inputs in unreduced form without incurring in extra costs, enabling the use of additions without correction or operations with incomplete reduction (see below for more details). We follow a similar procedure for computing field squaring. For ARM, we implement the integer multiplication using the schoolbook method. In this case, and also for modular reduction, we extensively exploit the parallelism of ARM and NEON instructions. The details are discussed in Section 4.4.

Let $0 \leq a, b < 2^m - c$. Field subtraction is computed as $(a - b) + borrow \cdot 2^m - borrow \cdot c$, where $borrow = 0$ if $a \geq b$, otherwise $borrow = 1$. Notice that in practice the addition with $borrow \cdot 2^m$ can be efficiently implemented by clearing the $(m + 1)$ -th bit of $a - b$.

Incomplete Reduction. Similar to [29], we exploit the form of the pseudo-Mersenne prime in combination with the incomplete reduction technique to speedup computations. We also mix incompletely reduced and completely reduced operands in novel ways.

Let $0 \leq a < 2^m - c$ and $0 \leq b < 2^m$. Field addition with incomplete reduction is computed as $(a + b) - carry \cdot 2^m + carry \cdot c$, where $carry = 0$ if $a + b < 2^m$, otherwise $carry = 1$. Again, in practice the subtraction with $carry \cdot 2^m$ can be efficiently implemented by clearing the $(m + 1)$ -th bit of $a + b$. The result is correct modulo p , but falls in the range $[0, 2^m - 1]$. Thus, this addition operation with incomplete reduction works with both operands in completely reduced form or with one operand in completely reduced form and one in incompletely reduced form. A similar observation applies to subtraction. Consider two operands a and b , such that $0 \leq a < 2^m$ and $0 \leq b < 2^m - c$. The standard field subtraction $(a - b) \pmod{2^m - c}$ described above will then produce an incompletely reduced result in the range $[0, 2^m - 1]$, since $a - b$ with $borrow = 0$ produces a result in the range $[0, 2^m - 1]$ and $a - b$ with $borrow = 1$ produces a result in the range $[-2^m + c + 1, -1]$, which is then fully reduced by adding $2^m - c$. Thus, performance can be improved by using incomplete reduction for an addition preceding a subtraction. For example, this technique is exploited in the point doubling computation (see Steps 7-8 of Algorithm 8). Note that, in contrast to addition, only the first operand is allowed to be in incompletely reduced form for subtraction.

To guarantee correctness in our software, and following the previous description, incompletely reduced results are always fed to one of the following: one of the operands of an incompletely reduced addition, the first operand of a field subtraction, a field multiplication or squaring (which ultimately produces a completely reduced output), or a field addition without correction preceding a field multiplication or squaring.

In the targeted setting, there are only a limited number of spots in the curve arithmetic in which incompletely reduced numbers cannot be efficiently exploited. For these few cases, we require a standard field addition. We use the efficient implementation proposed by Bos et al. [7]. Again, let $0 \leq a, b < 2^m - c$. Field addition is then computed as $((a + c) + b) - carry \cdot 2^m - (1 - carry) \cdot c$, where $carry = 0$ if $a + b + c < 2^m$, otherwise $carry = 1$. Similar to previous cases, the subtraction with $carry \cdot 2^m$ can be efficiently carried out by clearing the $(m + 1)$ -th bit in $(a + c) + b$. As discussed above, this efficient computation is also advantageously exploited in the modular reduction for multiplication and squaring.

4.3 Quadratic Extension Field Arithmetic

For the remainder, we use the following notation: (i) I, M, S, A and R represent inversion, multiplication, squaring, addition and modular reduction over \mathbb{F}_p , respectively, (ii) M_i and A_i represent integer multiplication and integer addition, respectively, and (iii) i, m, s, a and r represent analogous operations over \mathbb{F}_{p^2} . When representing registers in algorithms, capital letters are used to allocate operands with “double precision” (in our case, 256 bits). For simplification purposes, in the operation counting an integer operation with

double-precision is considered equivalent to two integer operations with single precision. We assume that addition, subtraction, multiplication by two and negation have roughly the same cost.

Let $a = a_0 + a_1i \in \mathbb{F}_{p^2}$ and $b = b_0 + b_1i \in \mathbb{F}_{p^2}$. Inversion over \mathbb{F}_{p^2} is computed as $a^{-1} = (a_0 - a_1i)/(a_0^2 + a_1^2)$. Addition and subtraction over \mathbb{F}_{p^2} consist in computing $(a_0 + b_0) + (a_1 + b_1)i$ and $(a_0 - b_0) + (a_1 - b_1)i$, respectively. We compute multiplication over \mathbb{F}_{p^2} using the Karatsuba method. In this case, we fully exploit lazy reduction and the room of one bit that is gained by using a prime of 127 bits. The details for the x64 implementation are shown in Algorithm 3. Remarkably, note that only the subtraction in Step 3 requires a correction to produce a positive result. No other addition or subtraction requires correction to positive or to modulo p . That is, \times , $+$ and $-$ represent operations over the integers. In addition, the algorithm accepts inputs in completely or incompletely reduced form and always produces a result in completely reduced form. Optionally, one may “delay” the computation of the final modular reductions (by setting *reduction* = *FALSE* in Alg. 3) if lazy reduction could be exploited in the curve arithmetic. This has been proven to be useful to formulas for the Weierstrass form [1], but unfortunately the technique cannot be advantageously exploited in the most efficient formulas for twisted Edwards (in this case, one should set *reduction* = *TRUE*). Squaring over \mathbb{F}_{p^2} is computed using the complex method. The details for the x64 implementation are shown in Algorithm 4. In this case, all the additions are computed as integer operations since, again, results can be let to grow up to 128 bits, letting subsequent multiplications take care of the reduction step.

Algorithm 3 Multiplication in \mathbb{F}_{p^2} with reduction ($m = 3M_i + 9A_i + 2R$) and without reduction ($m_u = 3M_i + 9A_i$), using completely or incompletely reduced inputs (x64 platform).

Input: $a = (a_0 + a_1i)$ and $b = (b_0 + b_1i) \in \mathbb{F}_{p^2}$, where $0 \leq a_0, a_1, b_0, b_1 \leq 2^{127} - 1, p = 2^{127} - c, c$ small.

Output: $a \cdot b \pmod{p} \in \mathbb{F}_{p^2}$.

1: $T_0 \leftarrow a_0 \times b_0$	$[0, 2^{254} >$
2: $T_1 \leftarrow a_1 \times b_1$	$[0, 2^{254} >$
3: $C_0 \leftarrow T_0 - T_1$	$< -2^{254}, 2^{254} >$
4: if $C_0 < 0$, then $C_0 \leftarrow C_0 + 2^{128} \cdot p$	$[0, 2^{255} >$
5: if <i>reduction</i> = <i>TRUE</i> , then $c_0 \leftarrow C_0 \pmod{p}$	$[0, p >$
6: $t_0 \leftarrow a_0 + a_1$	$[0, 2^{128} >$
7: $t_1 \leftarrow b_0 + b_1$	$[0, 2^{128} >$
8: $T_2 \leftarrow t_0 \times t_1$	$[0, 2^{256} >$
9: $T_2 \leftarrow T_2 - T_0$	$[0, 2^{256} >$
10: $C_1 \leftarrow T_2 - T_1$	$[0, 2^{256} >$
11: if <i>reduction</i> = <i>TRUE</i> , then $c_1 \leftarrow C_1 \pmod{p}$	$[0, p >$
12: return if <i>reduction</i> = <i>TRUE</i> then $a \cdot b = (c_0 + c_1i)$, else $a \cdot b = (C_0 + C_1i)$.	

Algorithm 4 Squaring in \mathbb{F}_{p^2} ($s = 2M + 1A + 2A_i$), using completely reduced inputs (x64 platform).

Input: $a = (a_0 + a_1i) \in \mathbb{F}_{p^2}$, where $0 \leq a_0, a_1 \leq p - 1, p = 2^{127} - c, c$ small.

Output: $a^2 \pmod{p} \in \mathbb{F}_{p^2}$.

1: $t_0 \leftarrow a_0 + a_1$	$[0, 2^{128} >$
2: $t_1 \leftarrow a_0 - a_1 \pmod{p}$	$[0, p >$
3: $c_0 \leftarrow t_0 \times t_1 \pmod{p}$	$[0, p >$
4: $t_0 \leftarrow a_0 + a_0$	$[0, 2^{128} >$
5: $c_1 \leftarrow t_0 \times a_1 \pmod{p}$	$[0, p >$
6: return $a^2 = (c_0 + c_1i)$.	

4.4 Extension Field Arithmetic on ARM: Efficient Interleaving of ARM-Based and NEON-Based Multiprecision Operations

Algorithm 5 Double 128-bit integer product with ARM and NEON interleaved (`double_mul_neonarm`).

Input: $a = \{a_i\}, b = \{b_i\}, c = \{c_i\}, d = \{d_i\}, i \in \{0, \dots, 3\}$.

Output: $(F, G) \leftarrow (a \times b, c \times d)$.

```

1:  $(F, G) \leftarrow (0, 0)$ 
2: for  $i = 0$  to  $1$  do
3:    $(C_0, C_1, C_2) \leftarrow (0, 0, 0)$ 
4:   for  $j = 0$  to  $3$  do
5:      $(C_0, F_{i+j}, C_1, F_{i+j+2}) \leftarrow (F_{i+j} + a_i b_j + C_0, F_{i+j+2} + a_{i+2} b_j + C_1)$            {done by NEON}
6:   for  $j = 0$  to  $3$  do
7:      $(C_2, G_{i+j}) \leftarrow G_{i+j} + c_j d_i + C_2$                                            {done by ARM}
8:    $(F_{i+4}, F_{i+6}, G_{i+4}) \leftarrow (F_{i+4} + C_0, C_1, C_2)$ 
9: for  $i = 2$  to  $3$  do
10:  for  $j = 0$  to  $3$  do
11:     $(C_2, G_{i+j}) \leftarrow G_{i+j} + c_j d_i + C_2$                                          {done by ARM}
12:   $G_{i+4} \leftarrow C_2$ 
13: return  $(F, G)$ 

```

Algorithm 6 Triple 128-bit integer product with ARM and NEON interleaved (`triple_mul_neonarm`).

Input: $a = \{a_i\}, b = \{b_i\}, c = \{c_i\}, d = \{d_i\}, e = \{e_i\}, f = \{f_i\}, i \in \{0, \dots, 3\}$.

Output: $(F, G, H) \leftarrow (a \times b, c \times d, e \times f)$.

```

1:  $(F, G, H) \leftarrow (0, 0, 0)$ 
2: for  $i = 0$  to  $3$  do
3:    $(C_0, C_1, C_2) \leftarrow (0, 0, 0)$ 
4:   for  $j = 0$  to  $3$  do
5:      $(C_0, F_{i+j}, C_1, G_{i+j}) \leftarrow (F_{i+j} + a_j b_i + C_0, G_{i+j} + c_j d_i + C_1)$            {done by NEON}
6:   for  $j = 0$  to  $3$  do
7:      $(C_2, H_{i+j}) \leftarrow H_{i+j} + e_j f_i + C_2$                                          {done by ARM}
8:    $(F_{i+4}, G_{i+4}, H_{i+4}) \leftarrow (C_0, C_1, C_2)$ 
9: return  $(F, G, H)$ 

```

Algorithm 7 Double modular reduction with ARM and NEON interleaved (`double_red_neonarm`).

Input: A prime $p = 2^{127} - c$, $a = \{a_i\}, b = \{b_i\}, i \in \{0, \dots, 7\}$.

Output: $(F, G) \leftarrow (a \bmod p, b \bmod p)$.

```

1:  $(F_i, G_i) \leftarrow (a_i, b_i)_{i \in \{0, \dots, 3\}}$ 
2:  $(C_0, C_1, C_2) \leftarrow (0, 0, 0)$ 
3: for  $j = 0$  to  $1$  do
4:    $(C_0, F_j, C_1, F_{j+2}) \leftarrow (F_j + a_{j+4} c + C_0, F_{j+2} + a_{j+6} c + C_1)$            {done by NEON}
5: for  $j = 0$  to  $3$  do
6:    $(C_2, G_j) \leftarrow G_j + b_{j+4} c + C_2$                                                {done by ARM}
7:  $(F_2, F_4, G_4) \leftarrow (F_2 + C_0, C_1, C_2)$ 
8:  $(F_0, G_0) \leftarrow (F_4 c + F_0, G_4 c + G_0)$ 
9: return  $(F, G)$ 

```

The potential performance gain when interleaving ARM and NEON operations is well-known. This feature was exploited in [6] to speed up the Salsa20 stream cipher. On the other hand, Sánchez and Rodríguez-

Henríquez [36] showed how to take advantage of NEON instructions to perform independent multiplications in operations over \mathbb{F}_{p^2} . In the following, we go a step further and show how to exploit the increasingly efficient capacity of modern ARM processors for executing ARM and NEON instructions “simultaneously” to implement multiprecision operations, such as multiplication, squaring and modular reduction, over \mathbb{F}_{p^2} . In other words, we exploit the fact that when ARM code produces a data hazard in the pipeline, the NEON unit may be ready to execute vector instructions, and vice versa. Note that loading/storing values from ARM to NEON registers still remains relatively expensive, so in order to achieve an effective performance improvement, one should carefully interleave *independent* operations while minimizing the loads and stores from one unit to the other. Hence, operations such as multiplication and squaring over \mathbb{F}_{p^2} are particularly friendly to this technique, given the availability of internal independent multiplications in their formulas. Thus, using this approach, we implemented:

- a double integer multiplier (`double_mul_neonarm`) detailed in Algorithm 5, which interleaves a single 128-bit multiplication using NEON and a single 128-bit multiplication using ARM,
- a triple integer multiplier (`triple_mul_neonarm`) detailed in Algorithm 6, which interleaves two single 128-bit multiplication using NEON and one single 128-bit multiplication using ARM, and
- a double reduction algorithm (`double_red_neonarm`) detailed in Algorithm 7, that interleaves a single modular reduction using NEON and a single modular reduction using ARM.

Note that integer multiplication is implemented using the schoolbook method, which requires one multiplication, two additions, one shift and one bit-wise AND per iteration. These operations were implemented using efficient fused instructions such as UMLAL, UMAAL, VMLAL and VSRA [28], which add the result of a multiplication or shift to the destination register in one single operation, reducing code size.

To validate the efficiency of our approach, we compared the interleaved algorithms above with standard implementations using only NEON or ARM. In all the cases, we observed a reduction of costs in favor of our novel interleaved ARM/NEON implementations (see Section 5 for benchmark results).

`Triple_mul_neonarm` is nicely adapted to the computation of multiplication over \mathbb{F}_{p^2} , since this operation requires three integer multiplications of 128 bits (Steps 1, 2 and 8 of Algorithm 3). For the case of squaring over \mathbb{F}_{p^2} , we use `double_mul_neonarm` to compute the two independent integer multiplications (Steps 3 and 5 of Algorithm 4). Finally, for each case we can efficiently use a `double_red_neonarm`. The final algorithms for ARM are shown as Algorithms 10 and 11 in Appendix B.

4.5 Point Arithmetic

In this section, we describe implementation details and our optimized formulas for the point arithmetic. We use as basis the most efficient set of formulas proposed by Hisil et al. [22], corresponding to the case $a = -1$, that uses a combination of homogeneous projective coordinates $(X : Y : Z)$ and extended homogeneous coordinates of the form $(X : Y : Z : T)$, where $T = XY/Z$.

The basic algorithms for computing point doubling and addition are shown in Algorithms 8 and 9, respectively. In these algorithms, we extensively exploit incomplete reduction (denoted by with \oplus, \ominus), following the details given in Section 4.2. To ease coupling of doubling and addition in the main loop of the scalar multiplication computation, we make use of Hamburg’s “extensible” strategy and output values $\{T_a, T_b\}$, where $T = T_a \cdot T_b$, at every point operation, so that a subsequent operation may compute coordinate T if required. Note that the cost of doubling is given by $4m + 3s + 5a$. We do not apply the usual transformation $2XY = (X + Y)^2 - (X^2 + Y^2)$ because in our case it is faster to compute one multiplication and one incomplete addition than one squaring, one subtraction and one addition. In the setting of variable-base scalar multiplication (see Alg. 2), the main loop of the evaluation stage consists of a doubling-addition computation, which corresponds to the successive execution of Algorithms 8 and 9. For this case, precomputed points are more efficiently represented as $(X + Y, Y - X, 2Z, 2T)$ (corresponding to setting `EXT_COORD = TRUE` in Alg. 9), so the cost of addition is given by $8m + 6a$.

Algorithm 8 Twisted Edwards point doubling over \mathbb{F}_{p^2} (DBL = $4m + 3s + 5a$).

Input: $P = (X_1, Y_1, Z_1)$.

Output: $2P = (X_2, Y_2, Z_2)$ and $\{T_a, T_b\}$ such that $T_2 = T_a \cdot T_b$.

1: $T_a \leftarrow X_1^2$	(X_1^2)
2: $t_1 \leftarrow Y_1^2$	(Y_1^2)
3: $T_b \leftarrow T_a \oplus t_1$	$(X_1^2 + Y_1^2)$
4: $T_a \leftarrow t_1 - T_a$	$(Y_1^2 - X_1^2)$
5: $Y_2 \leftarrow T_b \times T_a$	$(Y_2 = (X_1^2 + Y_1^2)(Y_1^2 - X_1^2))$
6: $t_1 \leftarrow Z_1^2$	(Z_1^2)
7: $t_1 \leftarrow t_1 \oplus t_1$	$(2Z_1^2)$
8: $t_1 \leftarrow t_1 \ominus T_a$	$(2Z_1^2 - (Y_1^2 - X_1^2))$
9: $Z_2 \leftarrow T_a \times t_1$	$(Z_2 = (Y_1^2 - X_1^2)[2Z_1^2 - (Y_1^2 - X_1^2)])$
10: $T_a \leftarrow X_1 \oplus X_1$	$(2X_1)$
11: $T_a \leftarrow T_a \times Y_1$	$(2X_1 Y_1)$
12: $X_2 \leftarrow T_a \times t_1$	$(X_2 = 2X_1 Y_1 [2Z_1^2 - (Y_1^2 - X_1^2)])$
13: return $2P = (X_2, Y_2, Z_2)$ and $\{T_a, T_b\}$ such that $T_2 = T_a \cdot T_b$.	

Algorithm 9 Twisted Edwards point addition over \mathbb{F}_{p^2} (ADD = $8m + 6a$, mADD = $7m + 7a$ or $8m + 10a$).

Input: $P = (X_1, Y_1, Z_1)$ and $\{T_a, T_b\}$ such that $T_1 = T_a \cdot T_b$. If $EXT_COORD = FALSE$ then $Q = (x_2, y_2)$, else $Q = (X_2 + Y_2, Y_2 - X_2, 2Z_2, 2T_2)$.

Output: $P + Q = (X_3, Y_3, Z_3)$ and $\{T_a, T_b\}$ such that $T_3 = T_a \cdot T_b$.

1: $T_1 \leftarrow T_a \times T_b$	(T_1)
2: if $EXT_COORD = FALSE$ then $T_2 = x_2 \oplus x_2, T_2 = T_2 \times y_2$	$(2T_2)$
3: $t_1 \leftarrow T_2 \times Z_1$	$(2T_2 Z_1)$
4: if $Z_2 = 1$ then $t_2 \leftarrow T_1 \oplus T_1$ else $t_2 \leftarrow T_1 \times 2Z_2$	$(2T_1 Z_2)$
5: $T_a \leftarrow t_2 - t_1$	$(T_a = \alpha = 2T_1 Z_2 - 2T_2 Z_1)$
6: $T_b \leftarrow t_1 \oplus t_2$	$(T_b = \theta = 2T_1 Z_2 + 2T_2 Z_1)$
7: $t_2 \leftarrow X_1 \oplus Y_1$	$(X_1 + Y_1)$
8: if $EXT_COORD = TRUE$ then $Y_3 = Y_2 - X_2$, else $Y_3 = y_2 - x_2$	$(Y_2 - X_2)$
9: $t_2 \leftarrow Y_3 \times t_2$	$(X_1 + Y_1)(Y_2 - X_2)$
10: $t_1 \leftarrow Y_1 - X_1$	$(Y_1 - X_1)$
11: if $EXT_COORD = TRUE$ then $X_3 = X_2 + Y_2$, else $X_3 = x_2 \oplus y_2$	$(X_2 + Y_2)$
12: $t_1 \leftarrow X_3 \times t_1$	$(X_2 + Y_2)(Y_1 - X_1)$
13: $Z_3 \leftarrow t_2 - t_1$	$\beta = (X_1 + Y_1)(Y_2 - X_2) - (X_2 + Y_2)(Y_1 - X_1)$
14: $t_1 \leftarrow t_1 \oplus t_2$	$\omega = (X_1 + Y_1)(Y_2 - X_2) + (X_2 + Y_2)(Y_1 - X_1)$
15: $X_3 \leftarrow T_b \times Z_3$	$(X_3 = \beta\theta)$
16: $Z_3 \leftarrow t_1 \times Z_3$	$(Z_3 = \beta\omega)$
17: $Y_3 \leftarrow T_a \times t_1$	$(Y_3 = \alpha\omega)$
18: return $P + Q = (X_3, Y_3, Z_3)$ and $\{T_a, T_b\}$ such that $T_3 = T_a \cdot T_b$.	

5 Performance Analysis and Experimental Results

In this section, we carry out the performance analysis of the proposed GLV-based scalar multiplication method using the GLV-SAC representation, and present benchmark results of our constant-time implementations of curve `Ted127-glv4` on x64 and ARM platforms. We also assess the performance improvement obtained with the proposed ARM/NEON interleaving technique. For our experiments, we targeted a 3.4GHz Intel Core i7-2600 Sandy Bridge processor and a 3.4GHz Intel Core i7-3770 Ivy Bridge processor, from the Intel family, and a Samsung Galaxy Note with a 1.4GHz Exynos 4 Cortex-A9 processor and an Arndale

Board with a 1.7GHz Exynos 5 Cortex-A15 processor, from the ARM family, both equipped with the NEON vector unit. The x64 implementation was compiled with Microsoft Visual Studio 2012 and ran on 64-bit Windows (Microsoft Windows 8 OS). In our experiments, we turned off Intel’s hyperthreading and Turbo Boost technologies; we averaged the cost of 10^4 operations which were measured with the timestamp counter instruction `rdtsc`. The ARM implementation was developed and compiled with the Android NDK (ndk8d) toolkit. In this case, we averaged the cost of 10^4 operations which were measured with the `clock_gettime()` function and scaled to clock cycles using the processor frequency.

First, we present timings for all the fundamental operations of scalar multiplication in Table 1. Implementation details for quadratic extension field operations and point operations over \mathbb{F}_{p^2} can be found in Section 4. “IR” stands for incomplete reduction and “extended” represents the use of the extended coordinates $(X + Y, Y - X, 2Z, 2T)$ to represent precomputed points. The four-dimensional decomposition of the scalar follows [30]. In particular, a scalar k is decomposed in smaller k_i s.t. $\max(|k_i|) < Cr^{1/4}$ for $0 \leq i \leq 3$, where r is the 251-bit prime order and $C = 179$ for our case (see §4.1). In practice, however, we have found that the bitlength of k_i is at most 63 bits for our targeted curve. The decomposition can be performed as a linear transformation by computing $k_i = \sum_{j=0}^3 \mathbf{round}(S_j k) \cdot M_{i,j}$ for $0 \leq i < 4$, where $M_{i,j}$ and S_j are integer constants. We truncate operands in the `round()` operation, adding enough precision to avoid loss of data. Thus, the computation involves a few multi-precision integer operations exhibiting constant-time execution.

Table 1. Cost (in cycles) of basic operations on curve `Ted127-glv4`.

Operation		ARM Cortex-A9	ARM Cortex-A15	Intel Sandy Bridge	Intel Ivy Bridge
\mathbb{F}_{p^2}	ADD with IR	20	19	12	12
	ADD	39	37	15	15
	SUB	39	37	12	12
	SQR	223	141	59	56
	MUL	339	185	78	75
	INV	13,390	9,675	6,060	5,890
ECC	DBL	2,202	1,295	545	525
	ADD	3,098	1,831	690	665
	mADD ($Z_1 = 1$)	2,943	1,687	622	606
	Φ endomorphism ($Z_1 = 1$)	3,118	1,724	745	712
	Ψ endomorphism ($Z_1 = 1$)	1,644	983	125	119
Misc	8-point LUT (extended)	291	179	83	79
	GLV-based LSB-set recoding	1,236	873	482	482
	4-GLV decomposition	756	430	305	290

Next, we analyze the cost of GLV-based variable-base scalar multiplication on curve `Ted127-glv4`. Based on Algorithm 2, this operation involves the computation of one Φ endomorphism, 2 Ψ endomorphisms, 3 additions and 4 mixed additions in the precomputation stage; 63 doublings, 63 additions, one mixed addition and 64 protected table lookups in the evaluation stage; and one inversion and 2 multiplications over \mathbb{F}_{p^2} for converting the final result to affine. In total, the cost is given by $1i + 833m + 191s + 769a + 64LUT8 + 4M + 9A$. This operation count does not include other additional computations, such as the recoding to the GLV-SAC representation or the decomposition to 4-GLV, which are relatively inexpensive (see Table 1).

Compared to [30], which uses a method based on the regular windowed recoding [33], the optimized GLV-SAC method for variable-base scalar multiplication allows us to save 181 multiplications, 26 squarings and 228 additions over \mathbb{F}_{p^2} . Additionally, it only requires 8 precomputed points, which involve 64 protected table lookups over 8 points (denoted by `LUT8`) during scalar multiplication, whereas the method in [30] requires 36 precomputed points, which involve 68 protected table lookups over 9 points. For example, this represents in practice a 17% speedup in the computation and a 78% reduction in the memory consumption of precomputation on curve `Ted127-glv4`.

Finally, in Table 2 we summarize our benchmark results for scalar multiplication and compare them with other constant-time implementations in the literature. The results for the representative variable-base scenario set a new speed record for protected curve-based scalar multiplication on x64 and ARM processors. In comparison with the previously fastest genus one implementation on x64 by Longa and Sica [30], which runs in 137,000 cycles, the presented result injects a cost reduction of 30% on a Sandy Bridge machine. Likewise, in comparison with the state-of-the-art genus 2 implementation by Bos et al. [7], our results are between 21%-24% faster on x64 processors. It is also between 17%-19% faster than the very recent implementation by Oliveira et al. [34] based on a binary GLS curve using the 2-GLV method⁴, and about 2 times faster than Bernstein et al.’s implementation using a Montgomery curve over \mathbb{F}_p [4]. Moreover, our results also demonstrate that the proposed techniques bring a dramatic reduction in the overhead for protecting against timing attacks. An unprotected version of our implementation computes a scalar multiplication in 87,000 cycles on the Sandy Bridge processor, which is only 9% faster than our protected version. In the case of ARM, our implementation of variable-base scalar multiplication on curve `Ted127-glv4` is 27% and 32% faster than Bernstein and Schwabe’s [6] and Hamburg’s [18] implementation (respect.) of `curve25519` on a Cortex-A9 processor. Note, however, that comparisons on ARM are particularly difficult. The implementation of [6] was originally optimized for Cortex-A8, and [18] does not exploit NEON.

To put our results in perspective, note that the original GLS paper [14] reported a scalar multiplication that ran in 0.76 the time of the best available implementation on x64 (Core 2 Duo) at that time, namely a Montgomery curve over \mathbb{F}_p [16]. However, the former implementation is not protected against timing attacks whereas the latter is protected. If, optimistically, one assumes a 10% overhead to protect [14], the ratio above would increase to at least 0.83. Our software, on the other hand, runs in only 0.63 and 0.49 the time of two contemporary implementations also based on the same Montgomery curve, namely [18] and [4], respectively, on another x64 processor (Sandy Bridge). Although a precise comparison is difficult (ratios are obtained on different x64 architectures, GLS implementation [14] and ours exploit different prime forms, have different endomorphism and precomputation costs, etc.) and part of the increase in the speedup can be attributed to moving from 2 to 4-GLV decomposition, there is a wide margin that makes clear the improvement obtained by using the proposed techniques. A similar experimental comparison for ARM is not available in the literature. To our knowledge, we report the first implementation of a GLV-based GLS curve on an ARM processor.

Finally, in our experiments to assess the improvement obtained with the proposed ARM/NEON interleaving technique on the Cortex-A9 processor, we observed speedups close to 17% and 24% in comparison with implementations exploiting only ARM or NEON instructions, respectively. Remarkably, for the same figures on the Cortex-A15, we observed speedups in the order of 34% and 35%, respectively. These experimental results confirm the significant performance improvement enabled by the proposed technique, which exploits the increasing capacity of the latest ARM processors for parallelizing ARM and NEON instructions.

Table 2. Cost (in 10^3 cycles) of implementations of variable-base scalar multiplication with protection against timing-type side-channel attacks at approximately 128-bit security level. Results are approximated to the nearest 10^3 cycles.

Work		ARM	ARM	Intel	Intel
Curve	Precomputations	Cortex-A9	Cortex-A15	Sandy Bridge	Ivy Bridge
<code>Ted127-glv4</code> (this work)	512 bytes (8 points)	417	244	96	92
<code>Ted127-glv4</code> , Longa-Sica [30]	2.25 KB (36 points)	-	-	137	-
Binary GLS $E/\mathbb{F}_{2^{254}}$, Oliveira et al. [34]	512 bytes (8 points)	-	-	115	113
Genus 2 Kummer C/\mathbb{F}_p , Bos et al. [7]	0	-	-	126 (*)	117
<code>Curve25519</code> , Bernstein et al. [4]	0	-	-	194 (*)	183 (*)
<code>Curve25519</code> , Bernstein et al. [6]	0	568 (*)	-	-	-
<code>Curve25519</code> , Hamburg [18]	0	616	-	153	-

(*) Source: eBACS [5].

⁴ In the case of *unprotected* software on x64, Oliveira et al. [34] hold the current speed record with 72,000 cycles on an Intel Sandy Bridge. Their protected version is significantly more costly and runs in about 115,000 cycles.

Acknowledgements: We would like to thank Joppe Bos, Craig Costello, Francisco Rodríguez-Henríquez and the reviewers for their useful comments that helped us improve the quality of this work. Also, we would like to thank Francisco Rodríguez-Henríquez for giving us access to the Arndale board for the development of the ARM implementation.

References

1. D.F. Aranha, K. Karabina, P. Longa, C. Gebotys, and J. López. Faster explicit formulas for computing pairings over ordinary curves. In K.G. Paterson, editor, *Advances in Cryptology - EUROCRYPT*, volume 6632, pages 48–68. Springer, 2011.
2. D. Bernstein. Cache-timing attacks on AES. 2005. <http://cr.yp.to/antiforgery/cachetiming-20050414.pdf>.
3. D. Bernstein, P. Birkner, M. Joye, T. Lange, and C. Peters. Twisted Edwards curves. In S. Vaudenay, editor, *Proceedings of Africacrypt 2008*, volume 5023 of *LNCS*, pages 389–405. Springer, 2008.
4. D. Bernstein, N. Duif, T. Lange, P. Schwabe, and B.-Y. Yang. High-speed high-security signatures. In B. Preneel and T. Takagi, editors, *Proceedings of CHES 2011*, volume 6917 of *LNCS*, pages 124–142. Springer, 2011.
5. D. Bernstein and T. Lange. eBACS: ECRYPT Benchmarking of Cryptographic Systems, accessed on Dec 12, 2013. <http://bench.cr.yp.to/results-dh.html>.
6. D. Bernstein and P. Schwabe. NEON crypto. In E. Prouff and P.R. Schaumont, editors, *Cryptographic Hardware and Embedded Systems - CHES 2012*, volume 7428 of *Lecture Notes in Computer Science*, pages 320–339. Springer, 2012.
7. J.W. Bos, C. Costello, H. Hisil, and K. Lauter. Fast cryptography in genus 2. In T. Johansson and P.Q. Nguyen, editors, *Advances in Cryptology - EUROCRYPT*, volume 7881 of *LNCS*, pages 194–210. Springer, 2013.
8. J.W. Bos, C. Costello, H. Hisil, and K. Lauter. High-performance scalar multiplication using 8-dimensional GLV/GLS decomposition. In G. Bertoni and J.-S. Coron, editors, *Cryptographic Hardware and Embedded Systems - CHES 2013*, volume 8086 of *LNCS*, pages 331–348. Springer, 2013.
9. D. Brumley and D. Boneh. Remote timing attacks are practical. In S. Mangard and F.-X. Standaert, editors, *Proceedings of the 12th USENIX Security Symposium*, volume 6225 of *LNCS*, pages 80–94. Springer, 2003.
10. A. Faz-Hernández, P. Longa, and A.H. Sánchez. Efficient and secure algorithms for GLV-based scalar multiplication and their implementation on GLV-GLS curves (extended version). In *Cryptology ePrint Archive, Report 2013/158*, 2013. Available at: <http://eprint.iacr.org/2013/158>.
11. M. Feng, B.B. Zhu, M. Xu, and S. Li. Efficient comb elliptic curve multiplication methods resistant to power analysis. In *Cryptology ePrint Archive, Report 2005/222*, 2005. Available at: <http://eprint.iacr.org/2005/222>.
12. M. Feng, B.B. Zhu, C. Zhao, and S. Li. Signed MSB-set comb method for elliptic curve point multiplication. In K. Chen, R. Deng, X. Lai, and J. Zhou, editors, *Proceedings of Information Security Practice and Experience (ISPEC 2006)*, volume 3903 of *LNCS*, pages 13–24. Springer, 2006.
13. S. D. Galbraith, X. Lin, and M. Scott. Endomorphisms for faster elliptic curve cryptography on a large class of curves. In *J. Cryptology*, volume 24(3), pages 446–469, 2011.
14. S.D. Galbraith, X. Lin, and M. Scott. Endomorphisms for faster elliptic curve cryptography on a large class of curves. In A. Joux, editor, *Advances in Cryptology - EUROCRYPT*, volume 5479 of *LNCS*, pages 518–535. Springer, 2009.
15. R.P. Gallant, J.L. Lambert, and S.A. Vanstone. Faster Point Multiplication on Elliptic Curves with Efficient Endomorphisms. In J. Kilian, editor, *Advances in Cryptology - CRYPTO*, volume 2139 of *LNCS*, pages 190–200. Springer, 2001.
16. Thomé E. Gaudry, P. The mpFq library and implementing curve-based key exchanges. In *SPEED 2007*, pages 49–64, 2007.
17. A. Guillevic and S. Ionica. Four dimensional GLV via the Weil restriction. volume 8269 of *LNCS*, pages 79–96. Springer, 2013.
18. M. Hamburg. Fast and compact elliptic-curve cryptography. In *Cryptology ePrint Archive, Report 2012/309*, 2012. Available at: <http://eprint.iacr.org/2012/309>.
19. D. Hankerson, K. Karabina, and A. Menezes. Analyzing the Galbraith-Lin-Scott point multiplication method for elliptic curves over binary fields. *IEEE Trans. Computers*, 58(10):1411–1420, 2009.
20. D. Hankerson, A. Menezes, and S. Vanstone. *Guide to elliptic curve cryptography*. Springer Verlag, 2004.
21. M. Hedabou, P. Pinel, and L. Beneteau. Countermeasures for preventing comb method against SCA attacks. In R. Deng, F. Bao, H. Pang, and J. Zhou, editors, *Proceedings of Information Security Practice and Experience (ISPEC 2005)*, volume 3439 of *LNCS*, pages 85–96. Springer, 2005.

22. H. Hisil, K. Wong, G. Carter, and E. Dawson. Twisted Edwards curves revisited. In J. Pieprzyk, editor, *Advances in Cryptology - ASIACRYPT*, volume 5350 of *LNCS*, pages 326–343. Springer, 2008.
23. Z. Hu, P. Longa, and M. Xu. Implementing 4-dimensional GLV method on GLS elliptic curves with j -invariant 0. *Designs, Codes and Cryptography*, 63(3):331–343, 2012. Available at: <http://eprint.iacr.org/2011/315>.
24. M. Joye and M. Tunstall. Exponent recoding and regular exponentiation algorithms. In M. Joye, editor, *Proceedings of Africacrypt 2003*, volume 5580 of *LNCS*, pages 334–349. Springer, 2009.
25. P. C. Kocher. Timing Attacks on Implementations of Diffie-Hellman, RSA, DSS, and Other Systems. In N. Kobitz, editor, *Advances in Cryptology - CRYPTO*, volume 1109 of *LNCS*, pages 104–113. Springer, 1996.
26. P. C. Kocher, J. Jaffe, and B. Jun. Differential Power Analysis. In M. Wiener, editor, *Advances in Cryptology - CRYPTO*, volume 1666 of *LNCS*, pages 388–397. Springer, 1999.
27. C.H. Lim and P.J. Lee. More flexible exponentiation with precomputation. In Y. Desmedt, editor, *Advances in Cryptology - CRYPTO*, volume 839 of *LNCS*, pages 95–107. Springer, 1994.
28. ARM Limited. ARM Architecture Reference Manual: ARMv7-A and ARMv7-R edition, 2012.
29. P. Longa and C. Gebotys. Efficient techniques for high-speed elliptic curve cryptography. In S. Mangard and F.-X. Standaert, editors, *Proceedings of CHES 2010*, volume 6225 of *LNCS*, pages 80–94. Springer, 2010.
30. P. Longa and F. Sica. Four-dimensional Gallant-Lambert-Vanstone scalar multiplication. In X. Wang and K. Sako, editors, *Advances in Cryptology - ASIACRYPT*, volume 7658 of *LNCS*, pages 718–739. Springer, 2012.
31. P. Longa and F. Sica. Four-dimensional Gallant-Lambert-Vanstone scalar multiplication. In *Journal of Cryptology (to appear)*, 2013.
32. B. Möller. Algorithms for multi-exponentiation. In S. Vaudenay and A.M. Youssef, editors, *Proceedings of SAC 2001*, volume 2259 of *LNCS*, pages 165–180. Springer, 2001.
33. K. Okeya and T. Takagi. The width- w NAF method provides small memory and fast elliptic curve scalars multiplications against side-channel attacks. In M. Joye, editor, *Proceedings of CT-RSA 2003*, volume 2612 of *LNCS*, pages 328–342. Springer, 2003.
34. T. Oliveira, J. López, D.F. Aranha, and F. Rodríguez-Henríquez. Lambda coordinates for binary elliptic curves. In G. Bertoni and J.-S. Coron, editors, *Cryptographic Hardware and Embedded Systems - CHES 2013*, volume 8086 of *LNCS*, pages 311–330. Springer, 2013.
35. D.A. Osvik, A. Shamir, and E. Tromer. Cache attacks and countermeasures: the case of AES. In D. Pointcheval, editor, *Topics in Cryptology - CT-RSA 2006*, volume 3860, pages 1–20. Springer, 2006.
36. A.H. Sánchez and F. Rodríguez-Henríquez. NEON implementation of an attribute-based encryption scheme. In M. Jacobson, M. Locasto, P. Mohassel, and R. Safavi-Naini, editors, *International Conference on Applied Cryptography and Network Security - ACNS 2013*, volume 7954 of *LNCS*, pages 322–338. Springer, 2013.
37. B. Smith. Families of fast elliptic curves from \mathbb{Q} -curves. In K. Sako and P. Sarkar, editors, *Advances in Cryptology - ASIACRYPT*, volume 8269 of *LNCS*, pages 61–78. Springer, 2013.
38. D. Weber and T.F. Denny. The solution of McCurley’s discrete log challenge. In H. Krawczyk, editor, *Advances in Cryptology - CRYPTO*, volume 1462 of *LNCS*, pages 458–471. Springer, 1998.
39. T. Yanik, E. Savaş, and Ç.K. Koç. Incomplete reduction in modular arithmetic. In *IEE Proc. of Computers and Digital Techniques*, volume 149(2), pages 46–52, 2002.
40. S.-M. Yen and M. Joye. Checking before output may not be enough against fault- based cryptanalysis. *IEEE Trans. Computers*, 49(9):967–970, 2000.
41. S.-M. Yen, S. Kim, S. Lim, and S.-J. Moon. A countermeasure against one physical cryptanalysis may benefit another attack. In K. Kim, editor, *Information Security and Cryptology - ICISC 2001*, volume 2288 of *Lecture Notes in Computer Science*, pages 414–427. Springer, 2002.

A Formulas for Endomorphisms Φ and Ψ on Curve Ted127-g1v4

Let $P = (X_1, Y_1, Z_1)$ be a point in homogeneous projective coordinates on a twisted Edwards curve with eq. (2), $u = 1 + i$ be a quadratic non-residue in \mathbb{F}_{p^2} , and $\zeta_8 = u/\sqrt{2}$ be a primitive 8th root of unity. Then, we can compute $\Phi(P) = (X_2, Y_2, Z_2, T_2)$ as follows

$$\begin{aligned} X_2 &= -X_1 (\alpha Y_1^2 + \theta Z_1^2) [\mu Y_1^2 - \phi Z_1^2], & Y_2 &= 2Y_1 Z_1^2 [\phi Y_1^2 + \gamma Z_1^2], \\ Z_2 &= 2Y_1 Z_1^2 [\mu Y_1^2 - \phi Z_1^2], & T_2 &= -X_1 (\alpha Y_1^2 + \theta Z_1^2) [\phi Y_1^2 + \gamma Z_1^2], \end{aligned}$$

where $\alpha = \zeta_8^3 + 2\zeta_8^2 + \zeta_8$, $\theta = \zeta_8^3 - 2\zeta_8^2 + \zeta_8$, $\mu = 2\zeta_8^3 + \zeta_8^2 - 1$, $\gamma = 2\zeta_8^3 - \zeta_8^2 + 1$ and $\phi = \zeta_8^2 - 1$.

For curve **Ted127-glv4**, we have the fixed values

$$\begin{aligned} \zeta_8 &= 1 + Ai, & \alpha &= A + 2i, & \theta &= A + Bi, \\ \mu &= (A - 1) + (A + 1)i, & \gamma &= (A + 1) + (A - 1)i, & \phi &= (B + 1) + i, \end{aligned}$$

where $A = 143485135153817520976780139629062568752$, $B = 170141183460469231731687303715884099729$.

Computing an endomorphism Φ with the formula above costs $12m + 2s + 5a$ or only $8m + 1s + 5a$ if $Z_1 = 1$. Similarly, we can compute $\Psi(P) = (X_2, Y_2, Z_2, T_2)$ as follows

$$X_2 = \zeta_8 X_1^p Y_1^p, \quad Y_2 = Z_1^{p^2}, \quad Z_2 = Y_1^p Z_1^p, \quad T_2 = \zeta_8 X_1^p Z_1^p.$$

Given the value for ζ_8 on curve **Ted127-glv4** computing an endomorphism Ψ with the formula above costs approximately $3m + 1s + 2M + 5A$ or only $1m + 2M + 4A$ if $Z_1 = 1$.

B Algorithms for Quadratic Extension Field Operations exploiting Interleaved ARM/NEON Multiprecision Operations

Below are the algorithms for multiplication and squaring over \mathbb{F}_{p^2} , with $p = 2^{127} - c$, for ARM platforms. They exploit functions interleaving ARM/NEON-based operations, namely `double_mul_neonarm`, `triple_mul_neonarm` and `double_red_neonarm`, detailed in Algorithms 5, 6 and 7, respectively.

Algorithm 10 Multiplication in \mathbb{F}_{p^2} using completely or incompletely reduced inputs, $m = 3M_i + 9A_i + 2R$ (ARM platform).

Input: $a = (a_0 + a_1i)$ and $b = (b_0 + b_1i) \in \mathbb{F}_{p^2}$, where $0 \leq a_0, a_1, b_0, b_1 \leq 2^{127} - 1, p = 2^{127} - c, c$ small.

Output: $a \cdot b \in \mathbb{F}_{p^2}$.

1: $t_0 \leftarrow a_0 + a_1$	[0, 2^{128} >
2: $t_1 \leftarrow b_0 + b_1$	[0, 2^{128} >
3: $(T_0, T_1, T_2) \leftarrow \text{triple_mul_neonarm}(a_0, b_0, a_1, b_1, t_0, t_1)$	[0, 2^{256} >
4: $C_0 \leftarrow T_0 - T_1$	< $-2^{254}, 2^{254}$ >
5: if $C_0 < 0$, then $C_0 \leftarrow C_0 + 2^{128} \cdot p$	[0, 2^{255} >
6: $T_2 \leftarrow T_2 - T_0$	[0, 2^{256} >
7: $C_1 \leftarrow T_2 - T_1$	[0, 2^{256} >
8: return $(c_0, c_1) \leftarrow \text{double_red_neonarm}(C_0, C_1)$	[0, p >

Algorithm 11 Squaring in \mathbb{F}_{p^2} using completely reduced inputs, $s = 2M + 1A + 2A_i$ (ARM platform).

Input: $a = (a_0 + a_1i) \in \mathbb{F}_{p^2}$, where $0 \leq a_0, a_1 \leq p - 1, p = 2^{127} - c, c$ small.

Output: $a^2 \in \mathbb{F}_{p^2}$.

1: $t_0 \leftarrow a_0 + a_1$	[0, 2^{128} >
2: $t_1 \leftarrow a_0 - a_1 \bmod p$	[0, p >
3: $t_2 \leftarrow a_0 + a_0$	[0, 2^{128} >
4: $(C_0, C_1) \leftarrow \text{double_mul_neonarm}(t_0, t_1, t_2, a_1)$	[0, p^2 >
5: return $a^2 = \text{double_red_neonarm}(C_0, C_1)$	[0, p >
