

2D-GLS: Faster and exception-free scalar multiplication in the GLS254 binary curve

Marius A. Aardal¹ and Diego F. Aranha¹

Department of Computer Science, Aarhus University, Denmark
{maardal,dfaranha}@cs.au.dk

Abstract. We revisit and improve performance of arithmetic in the binary GLS254 curve by introducing the 2D-GLS scalar multiplication algorithm. The algorithm includes theoretical and practice-oriented contributions of potential independent interest: (i) for the first time, a proof that the GLS scalar multiplication algorithm does not incur exceptions, such that faster incomplete formulas can be used; (ii) faster dedicated atomic formulas that alleviate the cost of precomputation; (iii) a table compression technique that reduces the storage needed for precomputed points; (iv) a refined constant-time scalar decomposition algorithm that is more robust to rounding. We also present the first GLS254 implementation for Armv8. With our contributions, we set new speed records for constant-time scalar multiplication by 6% and 34.5% on respectively 64-bit Intel and Arm platforms.

Keywords: Binary elliptic curves · Software implementation · GLS254.

1 Introduction

Elliptic Curve Cryptography (ECC) has become the *de facto* standard for instantiating public key cryptography, with security based on the conjectured-as-exponential hardness of solving discrete logarithms over elliptic curve groups (ECDLP problem). Scalar multiplication, in particular for the unknown point scenario, is the most expensive operation in cryptographic protocols with security guarantees based on the ECDLP. Since its introduction in 1985, there was plenty of research in finding efficient and secure implementation strategies, and choosing optimal parameters to improve performance [3,24,6].

An early milestone in this research was the idea due to Gallant-Lambert-Vanstone (GLV) [14] of exploiting efficient endomorphisms to accelerate scalar multiplication. In the large characteristic case with the curve $E : y^2 = x^3 + b$ defined over \mathbb{F}_p for prime p , it initially manifested as evaluating $\psi : (x, y) \rightarrow (\beta x, y)$ for β a non-trivial cube root of unity. The technique was later generalized to Galbraith–Lin–Scott (GLS) curves defined over \mathbb{F}_{p^2} , and to exploit two or more endomorphisms as in the FourQ curve [27] and the genus-2 case [5]. In the binary field case, Koblitz curves are the classical example of curves equipped with endomorphisms [32]; a class later extended to include binary GLS curves [17].

Beyond performance, implementation security is also an important research problem, especially on embedded targets where side-channel attacks are more

feasible. The classical countermeasure against side-channel attacks is to formulate the arithmetic in the most *regular* way possible, such as constant-time implementation against timing attacks. This translates to removing secret-dependent branches and memory accesses, and employing *complete* point addition formulas without corner cases [33]. In the case of exploiting endomorphisms, additional care needs to be taken for a correct and secure implementation, for example using the GLV-SAC recoding technique [11] or by explicitly proving correctness of the specific scalar multiplication algorithm [6,27].

In this paper, we revisit the implementation of scalar multiplication in binary GLS curves at the 128-bit security level by improving efficiency and correctness of implementations of the GLS254 curve. Our contributions are:

- A 2D variant of the GLS scalar multiplication algorithm that changes the computation/storage overhead. The 2D algorithm spends more precomputation to reduce the number of point additions in the main loop.
- The first proof of correctness for the GLS scalar multiplication algorithms that proves it to be exception-free. This means that there are no corner cases in the main loop of the algorithm (with exception of possibly the last iteration), which enables faster incomplete formulas.
- Faster dedicated formulas to reduce the cost of precomputation, and a table compression technique that exploits the endomorphism to reduce storage.
- A refined scalar decomposition algorithm that can be easily implemented in constant time. The algorithm has robust parity and length guarantees that fill some gaps in previous works [30].
- An efficient formulation of arithmetic in $\mathbb{F}_{2^{254}}$ targeting Arm processors. The field arithmetic uses the interleaved representation proposed in the CHES’16 Rump Session [31]. We take the opportunity to include this formulation in the formal research literature, initially presented informally. Furthermore, this also closes affirmatively a question posed in [22] about the efficiency of binary curves in Armv8 processors, deemed “unclear”.

With these contributions, we obtain speed records in the 64-bit Armv8 and Intel platforms, improving on previous results by 34.5% and 6%, respectively. While the latter speedup may seem small, we remark that it comes after decades of successful research in improving performance of ECC, so diminishing returns are expected. Due to the upcoming move to post-quantum cryptography, these techniques could be of limited practicality, but we believe they are relevant for applications not necessarily needing long-term security and involving the computation of many scalar multiplications, such as private set intersection protocols [34]. The proof and table compression may find further application in accelerating GLV/GLS scalar multiplication in pairing-based cryptography.

The rest of the document is organized as follows. Section 1 discusses preliminaries on binary GLS curves and their efficient implementation. Section 3 introduces the 2D-GLS algorithm and its correctness proof. Section 4 pushes these ideas further by presenting the scalar decomposition algorithm, followed by dedicated formulas in Section 5. Section 6 discusses the implementation of field arithmetic in Armv8, with experimental results in Section 7. The interested

reader will also find a treatment of point compression for binary GLS curves in the Appendix, together with formulas that did not fit in the main body.

2 Preliminaries

An ordinary *binary elliptic curve* in Weierstrass form is defined as

$$E/\mathbb{F}_q : y^2 + xy = x^3 + ax^2 + b \tag{1}$$

with $q = 2^m$ and coefficients $a, b \in \mathbb{F}_q, b \neq 0$. For any \mathbb{F}_{q^k} , the points $P = (x, y) \in \mathbb{F}_{q^k} \times \mathbb{F}_{q^k}$ that satisfy the equation form an abelian group $E_{a,b}(\mathbb{F}_{q^k})$ together with a point at infinity \mathcal{O} , which acts as the identity. The group law is denoted with additive notation $P + Q$, such that the scalar multiplication operation is written as kP .

2.1 Binary GLS curves

In the interest of defining notation for later use, we briefly summarize the theory of binary GLS curves from [17].

Let E be an ordinary binary curve as defined previously. From Hasse’s theorem, $\#E(\mathbb{F}_q) = q + 1 - t$ for some trace t satisfying $|t| \leq 2\sqrt{q}$. Pick some $a' \in \mathbb{F}_{q^2}$ with $\text{Tr}'(a') = 1$, where Tr' is the field trace from \mathbb{F}_{q^2} to \mathbb{F}_2 defined as $\text{Tr}'(c) = \sum_{i=0}^{2m-1} c^{2^i}$. It can be shown that $E' = E_{a',b}$ is the quadratic twist of E over \mathbb{F}_{q^2} with $\#E'(\mathbb{F}_{q^2}) = (q - 1)^2 + t^2$, and that E and E' are isomorphic over \mathbb{F}_{q^4} under an involutive twisting isomorphism ϕ .

An endomorphism ψ over \mathbb{F}_{q^2} can be constructed for E' by composing ϕ with the q -power Frobenius map π as $\psi = \phi\pi\phi^{-1}$. Evaluating ψ over points $P \in E'(\mathbb{F}_{q^2})$ only requires field additions [30].

E' would in this scenario be referred to as a binary GLS curve. If $\#E' = hr$ where h is a small cofactor and r is prime, it holds in the unique order- r subgroup $\mathcal{S} = E'(\mathbb{F}_{q^2})[r]$ that $\psi^2 = -1 \pmod r$. The map ψ restricted to \mathcal{S} has an eigenvalue μ such that for $P \in \mathcal{S}$, $\psi(P) = \mu P$.

2.2 λ -projective coordinates for GLS scalar multiplication

In [30], Oliveira et al. introduced the λ -projective coordinate system. To date, it is the most efficient point representation for binary elliptic curves. Given an affine point $P = (x, y)$ with $x \neq 0$, its λ -affine representation is (x, λ) with $\lambda = x + \frac{y}{x}$. The λ -projective representation of P is (X, Λ, Z) with $X = xZ$, $\Lambda = \lambda Z$ for some $Z \neq 0$. The point at infinity can now be represented as $\mathcal{O} = (1, 1, 0)$. The curve equation in (1) is correspondingly transformed to

$$(\Lambda^2 + \Lambda Z + aZ^2)X^2 = X^4 + bZ^4. \tag{2}$$

Algorithm 1: Constant-time scalar multiplication (Oliveira et al. [30])

Input : $P \in \mathcal{S}$ in λ -affine coordinates, $k \in [1, r - 1]$, window size w
Output: kP in λ -affine coordinates

- 1 Decompose k into subscalars k_1, k_2 .
- 2 $c_j \leftarrow 1 - (k_j \bmod 2)$ for $j = 1, 2$.
- 3 $k_j \leftarrow k_j + c_j$
- 4 Compute width- w length- ℓ odd signed regular recoding \bar{k}_1, \bar{k}_2 of k_1, k_2 .
- 5
- 6 Compute $T[i] = (2i + 1)P$ for all odd $i \in \{0, \dots, 2^{w-2} - 1\}$.
- 7 Convert T to λ -affine coordinates using a simultaneous inversion.
- 8
- 9 Perform a linear pass over T to recover $P_{j,\ell-1} = \bar{k}_{j,\ell-1}P$ for $j = 1, 2$.
- 10 $Q \leftarrow P_{1,\ell-1} + \psi(P_{2,\ell-1})$
- 11 **for** i **from** $\ell - 2$ **downto** 0 **do**
- 12 $Q \leftarrow 2^{w-2}Q$
- 13 Perform a linear pass over T to recover $P_{j,i} = \bar{k}_{j,i}P$ for $j = 1, 2$.
- 14 $Q \leftarrow 2Q + P_{1,i} + \psi(P_{2,i})$
- 15 $Q \leftarrow Q - c_1P - c_2\psi(P)$
- 16 Convert Q to λ -affine coordinates.
- 17 **return** Q ;

The constant-time binary GLS scalar multiplication algorithm from [30] is included in Algorithm 1. It is a constant-time left-to-right double-and-add algorithm using λ -projective coordinates, combining the Joye-Tunstall regular recoding algorithm [21] with the GLV interleaving technique. To use the GLV method, the scalar k is decomposed into two subscalars k_1, k_2 of roughly half length such that $k \equiv k_1 + k_2\mu \pmod{r}$. The two smaller scalar multiplications can then be computed in an interleaved fashion to save half of the point doublings.

In addition, a width- w windowing strategy is used. The subscalars are recoded into $\ell = \lceil m/(w - 1) \rceil$ odd signed digits of $w - 1$ bits using Algorithm 6 from [6] (which is a constant-length modification of Algorithm 6 from [21]). This means that both subscalars need to be made odd (line 3) and the result fixed in the end (line 16). By initially computing a table $T[i] = (2i + 1)P$ for all positive digits (in a phase known as the precomputation), the main loop can process the scalars one digit at a time, reducing the number of iterations by a factor $w - 1$. To be resistant against (cache-)timing attacks, each lookup requires a linear pass over the entire table, and there can be no branches dependent on c_j, k_j .

However, both [30] and the subsequent [31] suffer from a lack of rigor. First and foremost, no proof has been presented for correctness of the scalar multiplication algorithm. The λ -projective group law formulas are incomplete, so it might be that it can fail in corner cases. It also relies on *ad-hoc* tricks for constant-time scalar decomposition, with no proof of correctness or length guarantees.

2.3 GLS254 and the choice of parameters

Previous works have benchmarked their implementation of scalar multiplications over a GLS curve specially crafted for efficiency at the 128-bit security level. For the GLS254 curve, one chooses $m = 127$, such that the base field can be defined as $\mathbb{F}_q \equiv \mathbb{F}_2[z]/(z^{127} + z^{63} + 1)$ and its quadratic extension as $\mathbb{F}_{q^2} \equiv \mathbb{F}_q[u]/(u^2 + u + 1)$. The curve coefficients should be chosen to have minimal Hamming weight such that multiplying by them is as efficient as possible. We performed a parameter search that reproduced the curve chosen at [31]. By fixing $a' = u$ and searching for the shortest $b = (z^i + 1)$ such that the curve has order $2r$ for prime r , we were able to confirm that $i = 27$ is the smallest choice, giving a 254-bit r . This means that a multiplication by b can be computed with a single shifted addition.

To protect against Weil descent and generalized Gaudry-Hess-Smart (gGHS) attacks [15,19], several precautions must be taken. We pick m to be prime, as is the case for GLS254. In addition, the choice of b must be verified to not allow the attack, which happens with negligibly small probability for random b [17]. We used the MAGMA implementation of [8] available at <https://github.com/JJChiDguez/gGHS-check> to clear our particular choice. This particular check, together with the curve generation method geared towards efficiency, satisfies rigidity concerns [4]. We stress that the ECDLP in binary curves remains infeasible for the parameter range used in this work [12].

3 Scalar multiplication in GLS curves

In this section, we begin by presenting a new scalar multiplication variant for binary GLS curves. It combines the Shamir-Straus' trick [9] for multiple scalar multiplication with a new table compression technique using the GLS endomorphism. We refer to it as the 2D variant because it builds a two-dimensional table $T[i, j] = iP + j\psi(P)$ instead of $T[i] = iP$. Then, in subsection 3.2, we prove that the GLS scalar multiplication algorithms are exception-free.

3.1 The 2D variant

As in some *fast* variants of the Shamir-Straus' trick [9] for multiple scalar multiplication, the idea is to precompute $T[i, j] = iP + j\psi(P)$ for odd i, j . In the scalar multiplication loop, we then save roughly one point addition per iteration of the main loop by computing $2Q + T[i, j]$ instead of $2Q + T[i] + \psi(T[j])$.

This method was previously deemed noncompetitive due to the blowup in the size of the table. Because the subsingular regular recoding uses signed digits, we need to efficiently retrieve $s_1 iP + s_2 j\psi(P)$ for any $i, j \in \{1, \dots, 2^{w-1} - 1\}$ and sign combination $s_1, s_2 \in \{\pm 1\}$. The standard approach would be to build a table of $iP \pm j\psi(P)$ and then use conditional negations to get the two other combinations. The 2D table would then store $2^{2(w-2)+1}$ points. Even with specialized formulas for the precomputation, the cost in terms of storage and field operations is too high compared to the 1D algorithm.

The crucial new observation is that the efficiently computable GLS endomorphism ψ can also be used to compress the 2D table by a factor of 2. As $\psi^2(P) = -P$ for any $P \in \mathcal{S}$, we obtain the identity

$$-\psi(T[j, i]) = iP - j\psi(P).$$

It implies that we can generate all combinations from a table that only stores $iP + j\psi(P)$ for positive i, j . The rest of the combinations can be efficiently retrieved using conditional negations and conditional applications of ψ .

This compression trick not only halve the amount of precomputation needed, but also halves the time needed to do a linear pass through the table in the main loop. With new specialized group law formulas for the precomputation (see Section 5), the 2D algorithm is able to compete for the protected scalar multiplication speed record (see Section 7). The 2D variant is presented in Algorithm 2. For $w = 2$, the only difference is that a complete formula must be used for $2Q + P_1$ at $i = 1$ as well.

Algorithm 2: Constant-time 2D scalar multiplication

Input : $P \in \mathcal{S}$ in λ -affine coordinates, $k \in [1, r - 1]$, window size $w > 2$
Output: kP in λ -affine coordinates

- 1 Decompose k into odd subscalars k_1, k_2 using Algorithm 3.
- 2 Compute width- w length- ℓ odd signed regular recoding \bar{k}_1, \bar{k}_2 of k_1, k_2 . Here $\ell = \lceil \frac{w+1}{w-1} \rceil$.
- 3
- 4 Compute $T[i, j] = (2i + 1)P + (2j + 1)\psi(P)$ for all odd $i, j \in \{0, \dots, 2^{w-2} - 1\}$.
- 5 Convert T to λ -affine coordinates using a simultaneous inversion.
- 6
- 7 Perform a linear pass over T to recover $P_{\ell-1} = \bar{k}_{1, \ell-1}P + \bar{k}_{2, \ell-1}\psi(P)$
- 8 $Q \leftarrow P_{\ell-1}$
- 9 **for** i **from** $\ell - 2$ **downto** 1 **do**
- 10 $Q \leftarrow 2^{w-2}Q$
- 11 Perform a linear pass over T to recover $P_i = \bar{k}_{1, i}P + \bar{k}_{2, i}\psi(P)$
- 12 $Q \leftarrow 2Q + P_i$
- 13 Repeat the steps for $i = 0$, but use a complete formula for $2Q + P_0$.
- 14
- 15 Convert Q to λ -affine coordinates.
- 16 **return** Q ;

3.2 Proof of exception-free scalar multiplication

We will now prove that the scalar multiplication algorithms presented here and in [30] (with a minor modification) is correct on all valid inputs. The core issue is that the underlying λ -projective group law formulas from [30] are not complete, meaning that they output the wrong result in some corner cases. Without these exceptions, correctness would be trivial. One could instead explicitly handle these exceptions in constant time using complete formulas, but this would come at a high performance cost. Here, we prove that exceptional cases can only occur

in the last iteration(s) of the main loop. By using complete formulas there, correctness is ensured at only a minor performance penalty.

For clarity, the proof will be tailored to the 2D algorithm. However, it can be easily adapted to the 1D algorithm. The proof can be seen as a two-dimensional extension of the argument from Proposition 1 in [6]. We will for now assume that the scalar decomposition produces subscalars of bit-length at most $m + 1$, and defer the discussion about how to achieve this to Section 4. The proof crucially relies on the structure of the lattice discussed in [14,13] that emerge in the GLV method for scalar decomposition;

$$\mathcal{L} = \{(x, y) \in \mathbb{Z}^2 : x + y\mu \equiv 0 \pmod{r}\}.$$

Here r is the large prime order of $\#E'(\mathbb{F}_{q^2})$ and μ the eigenvalue of ψ restricted to \mathcal{S} . For our purposes, it is very useful to think about \mathcal{L} as the lattice of decompositions of zero (as done in [10]).

Lemma 1. *Let $q = 2^m$. Let E' be a binary GLS curve with $\#E'(\mathbb{F}_{q^2}) = hr$ for an odd prime r and small cofactor h . Let E be the curve defined over \mathbb{F}_q such that E' over \mathbb{F}_{q^2} is the quadratic twist of $E(\mathbb{F}_{q^2})$. Define*

$$v_1 = \left(\frac{(q-1)+t}{2}, \frac{(q-1)-t}{2} \right) \text{ and } v_2 = \left(\frac{(q-1)-t}{2}, \frac{-(q-1)-t}{2} \right),$$

where t is the trace of $E(\mathbb{F}_q)$. Then v_1, v_2 form an orthogonal basis for the lattice \mathcal{L} if and only if $h = 2$.

Proof. We first need to establish that $v_1, v_2 \in L$. As noted in [17], $\#E'(\mathbb{F}_{q^2}) = (q-1)^2 + t^2$. As the affine point $(0, \sqrt{b})$ of order 2 is always on the curve, so $2|h$. From $q = 2^m$ it now follows that t is odd. Then the numerators of the coordinates of v_1, v_2 are all even, meaning $v_1, v_2 \in \mathbb{Z}^2$. Using that $\mu \equiv (q-1)t^{-1} \pmod{r}$ (which follows from the same argument as in Lemma 2 of [13]),

$$\frac{(q-1)+t}{2} + \frac{(q-1)-t}{2} \mu \equiv \frac{(q-1)^2 + t^2}{2t} \equiv 0 \pmod{r}.$$

This shows that $v_1 \in \mathcal{L}$. Multiplying the lattice equation for v_1 by $-\mu$ and utilizing that $\mu^2 \equiv -1 \pmod{r}$, we get that $v_2 \in \mathcal{L}$ as well.

The next step is to show that v_1, v_2 form an orthogonal basis for \mathcal{L} , which per [13] has determinant r . It can be easily verified that v_1, v_2 are orthogonal vectors. This means that they form a basis of some lattice $\mathcal{L}' \subseteq \mathcal{L}$. The determinant of \mathcal{L}' is $\#E'(\mathbb{F}_{q^2})/2$. Hence, $\mathcal{L} = \mathcal{L}'$ if and only if $h = 2$. \square

We restrict the correctness proof to only the subclass of binary GLS curves with $h = 2$, which allows us to utilize the structure of \mathcal{L} revealed in Lemma 1. This is the only subclass one would care about in practice, as one wants \mathcal{S} to be as large as possible. We also require that m is prime, which is needed for security reasons anyways.

Theorem 1. *Let the notation be as in Lemma 1. Let $m > 4$ be prime and $2 \leq w \leq m$. Then Algorithm 2 is exception-free.*

Proof. Let us start by identifying the exceptional cases of the λ -projective formulas. The formula for $P+Q$ breaks down whenever $P = \pm Q$, $P = \mathcal{O}$ or $Q = \mathcal{O}$. \mathcal{O} does not have a λ -affine representation, so these last two cases are only a concern when the points are λ -projective. The $2P$ formula has no exceptional cases. Finally, the atomic formula for $2Q + P$ breaks down when $P = \pm 2Q$ or $Q = \mathcal{O}$.

We will argue that all the exceptions that can occur in Algorithm 1 encode an element of \mathcal{L} . By this, we mean that they define some $z_1, z_2 \in \mathbb{Z}$ such that $z_1P + z_2\psi(P) = \mathcal{O}$. Then $(z_1, z_2) \in \mathcal{L}$.

If $(z_1, z_2) \neq (0, 0)$, we can show that either $|z_1|$ or $|z_2|$ must be at least an m -bit integer. v_1, v_2 form an orthogonal basis of \mathcal{L} , and are both a solution to the SVP in \mathcal{L} with norm \sqrt{r} . Using the Hasse bound we get that $\|v_1\|, \|v_2\| \geq (q-1)/\sqrt{2}$. Now assume for contradiction that $|z_1|, |z_2| < q/2$. Then

$$\|(z_1, z_2)\| = \sqrt{z_1^2 + z_2^2} \leq \sqrt{2\left(\frac{q}{2} - 1\right)^2} = \frac{q-2}{\sqrt{2}} < \|v_1\|, \|v_2\|.$$

This is a contradiction, since v_1, v_2 have minimal norm in $\mathcal{L} - \{\mathbf{0}\}$. Thus, it must be the case that $|z_1| \geq q/2$ or $|z_2| \geq q/2$.

We now have all the tools needed to prove that no exceptions occur in Algorithm 2, and we will start with the precomputation stage. Assume that an exception did occur in the computation of $iP + j\psi(P)$ for some odd $i, j \in \{1, \dots, 2^{w-1} - 1\}$. $P, \psi(P) \neq \mathcal{O}$, so there can only be an exception if $iP = sj\psi(P)$ for some $s \in \{\pm 1\}$. Then $(i, -sj) \in \mathcal{L} - \{\mathbf{0}\}$. However, this is a contradiction, as neither $|i| = i$ nor $|-sj| = j$ are m -bit integers.

Next is the main loop. Let $Q_i = z_{1,i}P + z_{2,i}\psi(P)$ denote the value of Q after iteration i . No exception occurred in the precomputation stage, meaning Q is correctly initialized to $Q_{\ell-1} = \bar{k}_{1,\ell-1}P + \bar{k}_{2,\ell-1}\psi(P)$. Then at iteration i ,

$$Q_i = (2^{w-1}z_{1,i+1} + \bar{k}_{1,i})P + (2^{w-1}z_{2,i+1} + \bar{k}_{2,i})\psi(P).$$

Observe that as long as no exceptions occur, we have the invariant that

$$|z_{j,i+1}| \leq |z_{j,i}| \leq 2^{(\ell-i)(w-1)} - 1 \text{ and } z_{j,i} \neq 0 \text{ for } j = 1, 2.$$

Assume the first exception occurs at iteration i . The $w-2$ doublings are exception-free, so the exception must have been caused by the computation of $2Q_{i+1} + P_i$. Q_{i+1} cannot have been \mathcal{O} . This is because $2Q \neq \mathcal{O}$ for any $Q \neq \mathcal{O}$ and the incomplete $2Q+P$ formula does not output \mathcal{O} for any P, Q on the curve. Hence, the first exception must have occurred because $2^{w-1}Q_{i+1} = sP_i$ for an $s \in \{\pm 1\}$. This is equivalent to

$$(2^{w-1}z_{1,i+1} - s\bar{k}_{1,i})P + (2^{w-1}z_{2,i+1} - s\bar{k}_{2,i})\psi(P) = \mathcal{O}.$$

Define $z'_{j,i} = 2^{w-1}z_{j,i+1} - s\bar{k}_{j,i}$. Notice that $-s\bar{k}_{j,i}$ is a valid digit of the regular recoding. The invariants for $|z_{j,i}|$ then also hold for $|z'_{j,i}|$. Hence, $(z'_{1,i}, z'_{2,i}) \in \mathcal{L} - \{\mathbf{0}\}$.

Since m is prime, it holds for all $2 \leq w \leq m$ that

$$2^{(\lceil \frac{m}{w-1} \rceil - 1)(w-1)} - 1 \leq 2^{(\frac{m}{w-1} + \frac{w-2}{w-1} - 1)(w-1)} - 1 = 2^{m-1} - 1.$$

So it can't be the case that $i \geq \ell - \lceil \frac{m}{w-1} \rceil + 1$, because then $|z'_{1,i}|$ and $|z'_{2,i}|$ are of at most $m-1$ bits, which is a contradiction. Therefore, it must be the case that $i \leq \ell - \lceil \frac{m}{w-1} \rceil$.

For $w = 2$, this means that $i \leq 1$. For all other w , this means that $i = 0$. But these are exactly the iterations that use a complete formula for the computation of $2Q + P_i$ for the respective values of w . Thus, it is impossible for the first exception to occur in these last iterations. The conclusion is that there can be no exception in the main loop. \square

4 Scalar decomposition with parity & length guarantees

The GLV method for scalar decomposition needs a bit of care when required to run in constant-time while preserving length guarantees. The GLV method uses a reduced basis $\{u_1, u_2\}$ of some sublattice \mathcal{L}' of \mathcal{L} (see Section 3.2) to solve the CVP problem for $(k, 0)$ in \mathcal{L}' using Babai rounding [14]. For a given basis, there exist unique constants $N, \alpha_1, \alpha_2 \in \mathbb{Z}$ such that

$$(k, 0) = \beta_1 u_1 - \beta_2 u_2,$$

where $\beta_i = \frac{\alpha_i}{N} k$. The subscalars k_1, k_2 are then computed as

$$(k_1, k_2) = (k, 0) - b_1 u_1 - b_2 u_2,$$

where $b_i = \lceil \beta_i \rceil$. The magnitude of the subscalars can then be bounded by some expression depending on the norm of the basis vectors.

The issue for constant-time implementations is the computation of the b_i 's. Ideally we would compute them using divisions, but unfortunately divisions do not run in constant time in most processors ¹.

The standard solution was first introduced in [5] and further analyzed in [10]. The idea is to approximate the computation of the b_i 's using integer divisions by powers of 2, which can be implemented in constant-time using shifts. Choose some integer d such that $k < 2^d$, and precompute the constants $c_i = \lfloor \frac{\alpha_i}{N} 2^d \rfloor$. Then at runtime compute b_i as $b'_i = \lfloor \frac{c_i}{2^d} k \rfloor$.

This approach introduces rounding errors. As proven in Lemma 1 of [10], b'_i will either be b_i or incorrectly rounded down to $b_i - 1$. This does not affect the correctness of the decomposition. However, it does negatively impact the bounds on $|k_1|, |k_2|$. If the bounds become too loose, we might need more iterations of the main loop of the scalar multiplication to ensure correctness.

The standard sublattice basis used for scalar decomposition in binary GLS curves in the literature (e.g. in [2,30]) is from Lemma 3 of [13]. The basis from Lemma 1 is shorter by a factor of $\sqrt{2}$. If there were no rounding errors, we would

¹ See <https://www.bearssl.org/constanttime.html>.

get m -bit subscalars regardless of basis used. However, when taking rounding errors into account, our basis allows us to guarantee subscalars of $m + 1$ bits instead of $m + 2$. We prove this claim in Lemma 2.

That the rounding errors are one-sided can also be exploited to ensure that k_1, k_2 are odd, without affecting the length guarantees. This is desirable as the scalar-recoding algorithm only works for odd scalars. The alternative would be to initially set $k_i \leftarrow k_i + p_i$, where $p_i = 1 - (k_i \bmod 2)$. Then at the very end, one would compute $Q \leftarrow Q - p_1P - p_2\psi(P)$. With parities fixed, two point additions are saved and the scalar multiplication algorithm is simplified. Note that this holds for both the 1D and 2D scalar multiplication algorithms.

Algorithm 3: Constant-time odd scalar decomposition for binary GLS curves

Input : $k \in [1, r - 1]$
Constants: $N = \#E'(\mathbb{F}_{q^2})$, $\alpha_1 = q - 1 + t$, $\alpha_2 = q - 1 - t$
 $d = \lceil \frac{2m}{W} \rceil \cdot W$, where W is the size of a machine word.
 $c_i = \lfloor \frac{\alpha_i}{N} 2^d \rfloor$ for $i = 1, 2$
Output : Odd k_1, k_2 such that $k_1 + k_2\mu \equiv k \pmod{r}$

- 1 $b_i \leftarrow c_i k \ggg d$ for $i = 1, 2$.
- 2 $(k_1, k_2) \leftarrow (k, 0) - b_1v_1 - b_2v_2$
- 3 **if** $\alpha_1 \equiv 0 \pmod{4}$ **then**
- 4 | $(u_1, u_2) \leftarrow (v_2, v_1)$
- 5 **else**
- 6 | $(u_1, u_2) \leftarrow (v_1, v_2)$
- 7 $p_i \leftarrow k_i + 1 \pmod{2}$ for $i = 1, 2$
- 8 $(k_1, k_2) \leftarrow (k_1, k_2) - p_1u_1 - p_2u_2$
- 9 **return** k_1, k_2

Lemma 2. *Let the notation be as in Lemma 1 and assume that $h = 2$ and $m > 4$. Algorithm 3 on input $k \in [1, r - 1]$ outputs a valid decomposition k_1, k_2 . The subscalars are odd and $|k_1|, |k_2| < 2q$.*

Proof. Let k_1, k_2 denote the output of the GLV method on input k and let k'_1, k'_2 the output of Algorithm 3.

We start with correctness. Per definition, $(k, 0) + \mathcal{L} \in \mathbb{Z}^2/\mathcal{L}$ is the set of valid decompositions of k . Algorithm 2 produces (k'_1, k'_2) by adding integer multiples of v_1 and v_2 to $(k, 0)$. Hence, $(k'_1, k'_2) \in (k, 0) + \mathcal{L}$.

Next, let's bound the magnitude of the subscalars. The basis vectors are orthogonal with norm \sqrt{r} . By the same argument as in Lemma 3 of [13] it follows that $\|(k_1, k_2)\| \leq \sqrt{r}/2$. To make the analysis independent of r , we can upperbound it as $r \leq (q + 1)^2/2$ using the Hasse bound. Then $\|v_1\|, \|v_2\| \leq (q + 1)/\sqrt{2}$ and $\|(k_1, k_1)\| \leq (q + 1)/2$.

It can be easily verified that α_1, α_2, N are specified such that $(k, 0) = \beta_1v_1 + \beta_2v_2$. Since $k < r \leq (q + 1)^2/2 \leq q^2 \leq 2^d$, it follows from Lemma 1 of [10] that b'_i is either b_i or incorrectly rounded down to $b_i - 1$.

Let r_i be the bit that is 1 if such a rounding error occurred when computing b'_i . Let s_i be the bit that is 1 if v_i was subtracted from (k_1, k_2) at line 8. Then using the triangle inequality, we can derive the bound on the subscalars.

$$\begin{aligned}
 |k'_1|, |k'_2| &\leq \|(k'_1, k'_2)\| \\
 &= \left\| \sum_{i=1}^2 (\beta_i - (b_i - r_i + s_i))v_i \right\| \\
 &\leq \left\| \sum_{i=1}^2 (\beta_i - b_i)v_i \right\| + \|v_1\| + \|v_2\| \\
 &\leq \left(\frac{q+1}{2}\right) + 2\left(\frac{q+1}{\sqrt{2}}\right) \\
 &< 2q \qquad \qquad \qquad (\text{Assuming } m > 4)
 \end{aligned}$$

Finally, we will show that k'_1, k'_2 are odd. The proof of Lemma 1 establishes that t is odd. $\frac{(q-1)+t}{2} = \frac{(q-1)-t}{2} + t$, meaning exactly one of the the coordinates of v_1 are odd. By symmetry, only the other coordinate of v_2 is odd. Because $\alpha_1 = 2\left(\frac{(q-1)+t}{2}\right)$, $\alpha_1 \equiv 0 \pmod{4}$ exactly when the 1st coordinate of v_2 is odd. Then u_i is the basis vector with the odd i -th coordinate. Subtracting (k'_1, k'_2) by u_i flips the parity of k'_i but leaves the parity of the other subscalar unchanged. $p_i = 1$ exactly when k_i is even, meaning that the subscalars output are odd. \square

Corollary 1. *Let $m > 4$ be a prime number. For any window size $2 < w \leq m$, the number of digits needed to recode the subscalars output by Algorithm 3 is the same as one would need for the subscalars output by the GLV method with no rounding errors. For $w = 2$, one more digit is required.*

5 New formulas for faster precomputation

The 2D scalar multiplication variant represents a different strategy for utilizing precomputation: the table grows quadratically faster than its 1D counterpart, which puts a lot more importance on reducing the cost of precomputation. For both the 1D and 2D variant, we present much more efficient strategies for the precomputation stage. The 2D precomputation (Algorithm 5) uses the 1D precomputation 4) as a subroutine, which makes a case for the fairness of our optimization efforts.

The precomputation algorithms depend on several new atomic group law formulas. Compared to doing the operations using the existing formulas from [30], they provide a significant saving in the number of field multiplications and squarings needed. Because these formulas are derived by combining the original group law formulas, they do not introduce additional exceptions. The new formulas that are nontrivial to derive are included in Appendix B.

It follows from the same argument as in Theorem 1 that the new precomputation algorithms are exception-free. At any step we compute $iP + j\psi(P)$ for

small coefficients i, j , where at least one of them are nonzero. Then $(i, \pm j) \notin \mathcal{L}$, meaning that there can be no exception.

Algorithm 4: Precomputation-1D

Input : $P \in \mathcal{S}$ in λ -affine coordinates, window size $w > 2$.
Output: Table T of size $2^{w-2} \times 2^{w-2}$ with $T[i] = (2i + 1)P$ in λ -affine coordinates.

- 1 $T[0] \leftarrow P$
- 2 $T[1] \leftarrow 3P$
- 3 **for** i **from** 0 **to** $2^{w-3} - 2$ **do**
- 4 $T[2i + 3], T[2i + 2] \leftarrow 2T[i + 1] \pm P$
- 5 Convert T to λ -affine coordinates using simultaneous inversion.
- 6 **return** T

Algorithm 5: Precomputation-2D

Input : $P \in \mathcal{S}$ in λ -affine coordinates, window size $w > 2$.
Output: Table T of size $2^{w-2} \times 2^{w-2}$ with $T[i, j] = (2i + 1)P + (2j + 1)\psi(P)$ in λ -affine coordinates.

- 1 $R \leftarrow \text{Precomputation-1D}(P, w)$
- 2 **for** i **from** 0 **to** $2^{w-2} - 1$ **do**
- 3 $T[i, i] \leftarrow R[i] + \psi(R[i])$
- 4 **for** j **from** 1 **to** $2^{w-2} - 1$ **do**
- 5 $Q \leftarrow \psi(R[j])$
- 6 **for** i **from** 0 **to** $j - 1$ **do**
- 7 $T[i, j], T[j, i] \leftarrow R[i] \pm Q$
- 8 $T[j, i] \leftarrow \psi(T[j, i])$
- 9 Convert T to λ -affine coordinates using simultaneous inversion.
- 10 **return** T

6 Binary field arithmetic for Arm

This section details our Arm implementation of the GLS254 curve. The focus will be on the field arithmetic. It was implemented specifically for the platform, relying heavily on 128-bit Arm Neon vector instructions to achieve high performance. The rest of the curve implementation is almost exclusively written in C, and therefore does not differ much from the Intel implementation from [30].

Specifically, our implementation targets Armv8 AArch64, which introduces some new useful instructions for cryptographic implementations. In particular, we take advantage of the new PMULL vector instruction for 64-bit binary polynomial multiplication, a direct analogue of PCLMULQDQ for Intel. For convenience of implementation, we use C intrinsics for the Arm Neon vector instructions.

Table 1. Cost of the λ -projective group law formulas with respect to the number of multiplications, multiplications by curve coefficients a and b and squarings $(\tilde{m}, \tilde{m}_a, \tilde{m}_b, \tilde{s})$ over the extension field \mathbb{F}_{q^2} . For the mixed point representations, Q is λ -projective while P , P_1 and P_2 are λ -affine. The formulas that have not been derived or that provided insignificant speedups are marked with '-'.

Op. \ Rep.	Projective	Mixed	Affine
$2P$	$4\tilde{m} + \tilde{m}_a + 4\tilde{s}/3\tilde{m} + 4\tilde{m}_a + \tilde{m}_b + 4\tilde{s}$	-	$\tilde{m} + 3\tilde{s}$
$3P$	-	-	$4\tilde{m} + \tilde{m}_a + 4\tilde{s}$
$P + Q$	$11\tilde{m} + 2\tilde{s}$	$8\tilde{m} + 2\tilde{s}$	$5\tilde{m} + 2\tilde{s}$
$P \pm Q$	-	$12\tilde{m} + 5\tilde{s}$	$6\tilde{m} + 4\tilde{s}$
$2Q+P$	-	$10\tilde{m} + \tilde{m}_a + 6\tilde{s}$	-
$2Q + P_1 + P_2$	-	$17\tilde{m} + \tilde{m}_a + 8\tilde{s}$	-
$P + \psi(P)$	-	-	$3.5\tilde{m} + 1.5\tilde{s}$

This section first details the implementation of the base field \mathbb{F}_q with $q = 2^m$ and $m = 127$, then how we implement the quadratic extension field \mathbb{F}_{q^2} on top.

6.1 Arithmetic in the base field \mathbb{F}_q

Representation of elements. One benefit of the choice of field, is that the bit vector representation of $a \in \mathbb{F}_q$ is 127 bits long, meaning that we can fit it in a single 128-bit Neon vector register. We denote $a[0], a[1]$ as respectively the least significant and most significant word of the 128-bit register that stores $a \in \mathbb{F}_{2^{127}}$. $a[0]$ stores the bit vector for the terms z^0 to z^{63} , $a[1]$ terms z^{64} to z^{126} . We will sometimes use the notation $a = \{a[0], a[1]\}$ to show the contents of the register.

An efficiency issue for Arm AArch64, compared to AArch32, is that it cannot reference the upper word $a[1]$ of a 128-bit register as a separate 64-bit register [16]. Instead, one needs to use the Arm Neon instruction `EXT`. It takes two registers a, b and outputs $\{a[1], b[0]\}$. The lower half of this output can then be referenced for further computation. Table 2 gives an overview of all the Neon instructions that we used for our implementation.

Table 2. Arm Neon 128-bit vector instructions used. The first 128-bit operand is denoted a , the second b . The output is also stored in a 128-bit register.

Symbol	Description	Neon Instruction
\oplus, \wedge	Bitwise XOR, AND	<code>EOR, AND</code>
\ll_{128}, \gg_{128}	Logical shift (no carry between words)	<code>SHL, SHR</code>
<code>pmull.bot</code>	Multiply binary polynomials $a[0], b[0]$	<code>PMULL</code>
<code>pmull.top</code>	Multiply binary polynomials $a[1], b[1]$	<code>PMULL2</code>
<code>extract</code>	Outputs $\{a[1], b[0]\}$	<code>EXT</code>

Polynomial multiplication. The polynomial multiplication algorithm takes as input two binary polynomials $a, b \in \mathbb{F}_q$ and outputs their polynomial product c . The degree of c can be up to twice the degree of the operands. Hence it must be stored in two 128-bit vector registers c_0, c_1 , where c_0 stores the lower half.

For polynomial multiplication we use the Arm Neon implementation from [16]. It efficiently performs 128-bit polynomial multiplication using the new PMULL instructions. While they managed to implement it using 3 multiplications with the Karatsuba algorithm on AArch32, the high number of EXTs this would require on AArch64 meant that they instead opted for an algorithm with an extra PMULL.

Polynomial squaring and field multi-squaring. Polynomial squaring of an $a \in \mathbb{F}_q$ can be trivially implemented as $c_0 \leftarrow \text{pmull_bot}(a, a)$, $c_1 \leftarrow \text{pmull_top}(a, a)$.

For multi-squaring in settings where it does not need to be computed in constant time, we implemented the technique from [1,7]. It uses lookup tables that are precomputed offline to compute the reduced result of a^{2^k} . However, for smaller Arm processors like the Cortex-A55, this method only outperforms the naive loop implementation for $k > 12$. This is a lot higher than the threshold of $k > 5$ for Intel [30]. It is an example of the higher cost of memory access on smaller Arm devices, which often result in a lower yield for precomputation based techniques.

Modular reduction. To compute a field multiplication or squaring, the result of the polynomial algorithm must be reduced modulo $f(z)$. We here present novel algorithms for efficient modular reduction, using exclusively Arm Neon vector instructions. Like the polynomial multiplication algorithm from [16], they attempt to minimize the number of accesses to the top half of the 128-bit registers, each of which incurs the cost of an EXT.

For reducing the result of a polynomial multiplication, we use Algorithm 6. It implements the lazy reduction technique from [31]. Instead of reducing $f(z)$, we reduce by the redundant trinomial $z \cdot f(z) = z^{128} + z^{64} + z$. Reductions by $zf(z)$ are roughly 40% faster than proper reductions by $f(z)$. The result can have degree up to 127 instead of 126, but as the result still fits in a 128-bit register, this makes no difference. As $(c \bmod zf(z)) \bmod f(z) = c \bmod f(z)$, one can easily recover the properly reduced result from the output of the lazy reduction. This is done by conditionally adding $z^{63} + 1$ to it when bit 127 is set.

It is possible to reduce the product of a squaring slightly faster using Algorithm 7. It exploits the fact that after a polynomial squaring, c_0, c_1 only have bits set at even positions. Thus, we can remove the logic from Algorithm 6 that handles the carry of bit 191 for the left shift by 1, since bit 191 is always 0. It can also be observed that as bit 127 and 191 are 0, Algorithm 7 actually computes the proper reduction modulo $f(z)$.

Field inversion. Field inversion is done in the same way as for the Intel implementation in [30], using the Itoh-Tsujii algorithm [20]. We generated our addition

Algorithm 6: Lazy reduction by $z \cdot f(z) = z^{128} + z^{64} + z$

Input : 254-bit polynomial stored in two 128-bit registers c_0, c_1 .

Output: 128-bit register a storing $c(z) \bmod f(z)$.

Temps.: Uses 128-bit registers t_0, t_1, t_2 .

```

1  $t_0[0] \leftarrow 0$ 
2  $t_1[0] \leftarrow c_1[0] \gg 63$ 
3  $t_0 \leftarrow \text{extract}(c_1, t_0)$ 
4  $t_2[0] \leftarrow c_1[0] \oplus t_0[0]$ 
5  $t_1[0] \leftarrow t_1[0] \oplus t_2[0]$ 
6  $t_0 \leftarrow \text{extract}(t_0, t_1)$ 
7  $a \leftarrow c_0 \oplus t_0$ 
8  $t_2 \leftarrow t_2 \ll_{128} 1$ 
9  $a \leftarrow a \oplus t_2$ 
10 return  $a$ 

```

Algorithm 7: Reduction by $f(z) = z^{127} + z^{63} + 1$ in \mathbb{F}_q after squaring

Input : 253-bit polynomial c stored in two 128-bit registers c_0, c_1 .

Output: 128-bit register a storing $c(z) \bmod f(z)$.

Temps.: Uses 128-bit registers t_0, t_1 .

```

1  $t_0[0] \leftarrow 0$ 
2  $t_0 \leftarrow \text{extract}(c_1, t_0)$ 
3  $t_1[0] \leftarrow c_1[0] \oplus t_0[0]$ 
4  $t_0 \leftarrow \text{extract}(t_0, t_1)$ 
5  $a \leftarrow c_0 \oplus t_0$ 
6  $t_1 \leftarrow t_1 \ll_{128} 1$ 
7  $a \leftarrow a \oplus t_1$ 
8 return  $a$ 

```

chain for $m - 1 = 126$ using McLoughlin's *addchain* library [29].

$$1 \rightarrow 2 \rightarrow 3 \rightarrow 6 \rightarrow 12 \rightarrow 24 \rightarrow 30 \rightarrow 48 \rightarrow 96 \rightarrow 126$$

The cost of a field inversion is therefore $m - 1$ squarings and 9 multiplications. The steps after 30 involve multi-squarings with $k > 12$. When the inversion does not have to be in constant time, these steps can then be sped up using the table-based multi-squaring approach.

6.2 Arithmetic in the extension field \mathbb{F}_{q^2}

The elements of this field can be represented as polynomials $a_1u + a_0$, with coefficients $a_0, a_1 \in \mathbb{F}_q$. Therefore, we need two 128-bit registers to represent them. The extension field arithmetic can be implemented from the base field arithmetic, using the identities presented in [30].

In [31], Oliveira et al present an algorithm for simultaneously reducing both coefficients of an element in \mathbb{F}_{q^2} at the cost of only a single base field reduction. We have included the Arm Neon implementation in Algorithm 8.

Algorithm 8: Lazy simultaneous reduction by $zf(z) = z^{128} + z^{64} + z$ for coordinate-wise reduction in \mathbb{F}_{q^2} (Oliveira et al. [31])

Input : Unreduced polynomial stored in interleaved 128-bit registers c_0, c_1, c_2, c_3 .
Output: 128-bit register a storing $c(z) \bmod zf(z)$.
Temps.: Uses 128-bit register t .

- 1 $c_2 \leftarrow c_2 \oplus c_3$
- 2 $t \leftarrow c_3 \ll_{128} 1$
- 3 $c_1 \leftarrow c_1 \oplus t$
- 4 $c_1 \leftarrow c_1 \oplus c_2$
- 5 $t \leftarrow c_2 \gg_{128} 63$
- 6 $c_1 \leftarrow c_1 \oplus t$
- 7 $t \leftarrow t \ll_{128} c_2$
- 8 $c_0 \leftarrow c_0 \oplus t$
- 9 **return** c_0, c_1

However, the reduction algorithm requires the field elements to be stored in an interleaved representation. For an $a \in \mathbb{F}_{q^2}$, let a_0, a_1 be the 128-bit registers storing each of its coefficients. Then the interleaved representation of a is

$$a'_0 = \{a_0[0], a_1[0]\}, \quad a'_1 = \{a_0[1], a_1[1]\}. \quad (3)$$

Note that a'_0, a'_1 store a_0 in the lower half and a_1 in the upper half. The input to the reduction algorithm must also be in an interleaved representation. Let c_0, c_1, c_2, c_3 be the non-interleaved 128-bit output registers of a polynomial multiplication or squaring computed using the identities from [30]. Then c_0, c_1 store the unreduced constant coefficient and c_2, c_3 the other. The interleaved representation c'_0, c'_1, c'_2, c'_3 of these unreduced coefficients continue the pattern from (3). c'_0, c'_1 are c_0, c_2 interleaved, and c'_2, c'_3 are c_1, c_3 interleaved.

In order to reap the benefits of the reduction algorithm, we had to implement the \mathbb{F}_{q^2} arithmetic directly in the interleaved representation. To do this, we manually merged and interleaved the base field algorithms to compute the identities from [30]. The only exception is inversion, where a standard base field inversion is used as a subroutine, which again uses all the arithmetic operations discussed in the previous section. While the abstraction between base field and extension field is somewhat broken for the sake of performance, the base field implementation is still the crucial foundation.

7 Results and discussion

Our implementations, together with SAGE scripts for verification and operation counts, can be found at <https://github.com/dfaranha/gls254>.

Table 3. The cost of the scalar multiplications with respect to the number of inversions, multiplications and squarings $(\tilde{i}, \tilde{m}, \tilde{s})$ over \mathbb{F}_{q^2} . The total cost in field multiplications are estimated using $\tilde{i}_{\text{non-ct}} = 18\tilde{m}$, $\tilde{i}_{\text{ct}} = 27\tilde{m}$ and $\tilde{s} = 0.4\tilde{m}$, as measured in our platforms, and rounded to the nearest integer. The label “prev” denotes operation counts for previous work.

Variant \ w	3	4	5	6
Precomp.				
1D (prev)	$\tilde{i} + 12\tilde{m} + 6\tilde{s}$	$\tilde{i} + 38\tilde{m} + 14\tilde{s}$	$\tilde{i} + 90\tilde{m} + 30\tilde{s}$	$\tilde{i} + 194\tilde{m} + 62\tilde{s}$
1D	$\tilde{i} + 6\tilde{m} + 4\tilde{s}$	$\tilde{i} + 31\tilde{m} + 13\tilde{s}$	$\tilde{i} + 81\tilde{m} + 31\tilde{s}$	$\tilde{i} + 181\tilde{m} + 67\tilde{s}$
2D	$2\tilde{i} + 36\tilde{m} + 11\tilde{s}$	$2\tilde{i} + 158\tilde{m} + 43\tilde{s}$	$2\tilde{i} + 594\tilde{m} + 155\tilde{s}$	$2\tilde{i} + 2234\tilde{m} + 571\tilde{s}$
Main loop				
1D (both)	$1273\tilde{m} + 764\tilde{s}$	$979\tilde{m} + 680\tilde{s}$	$819\tilde{m} + 628\tilde{s}$	$738\tilde{m} + 608\tilde{s}$
2D	$823\tilde{m} + 633\tilde{s}$	$676\tilde{m} + 591\tilde{s}$	$593\tilde{m} + 561\tilde{s}$	$554\tilde{m} + 553\tilde{s}$
Total				
1D (prev)	$2\tilde{i} + 1309\tilde{m} + 780\tilde{s}$	$2\tilde{i} + 1041\tilde{m} + 704\tilde{s}$	$2\tilde{i} + 933\tilde{m} + 668\tilde{s}$	$2\tilde{i} + 956\tilde{m} + 680\tilde{s}$
1D	$2\tilde{i} + 1281\tilde{m} + 768\tilde{s}$	$2\tilde{i} + 1012\tilde{m} + 693\tilde{s}$	$2\tilde{i} + 902\tilde{m} + 659\tilde{s}$	$2\tilde{i} + 921\tilde{m} + 675\tilde{s}$
2D	$3\tilde{i} + 861\tilde{m} + 644\tilde{s}$	$3\tilde{i} + 839\tilde{m} + 634\tilde{s}$	$3\tilde{i} + 1189\tilde{m} + 716\tilde{s}$	$3\tilde{i} + 2790\tilde{m} + 1124\tilde{s}$
Est. mult.				
1D (prev)	$1666\tilde{m}$	$1368\tilde{m}$	$1245\tilde{m}$	$1273\tilde{m}$
1D	$1633\tilde{m}$	$1334\tilde{m}$	$1211\tilde{m}$	$1236\tilde{m}$
2D	$1182\tilde{m}$	$1155\tilde{m}$	$1538\tilde{m}$	$3303\tilde{m}$

7.1 Operation counts for binary GLS scalar multiplication

Throughout our work, we have used field operation counts as a measure of complexity of the scalar multiplication variants. With an understanding of the relative cost of the operations, the count gives a platform-independent estimate of the relative performance of the algorithms. In particular, it guided our choice of window size. However, it crucially does not capture the space-time trade-offs of a particular architecture. For the variants discussed in this paper, which precisely are variations in how to use space, this trade-off has a huge impact. This will be apparent in the next subsection.

Table 3 gives an overview the operation counts for the variants. Additions and multiplications by the curve coefficients are ignored due to their insignificant impact on performance. We include the costs of the 1D algorithm without the new formulas and scalar decomposition to highlight the impact of our contributions. For the sake of fairness, it has been modified to be exception-free in the same way as the others.

As expected, the 2D algorithm spends more time on precomputation and less in the main loop. We see that the model predicts $w = 5$ to be the sweet spot for 1D and $w = 4$ for 2D. Notably, 2D $w = 3$ is predicted to be faster than 1D $w = 5$ while using only half the space.

For a simpler comparison, we estimate the total cost in terms of field multiplications. The relative costs of the other operations are very close to what

we have benchmarked on all platforms. With this, the model predicts that 1D with $w = 5$ should be 2.7% faster from our contributions. The 2D approach with $w = 4$ is predicted to be 7.2% faster than 1D in previous work with $w = 5$, and 4.6% faster than 1D with $w = 5$ from this work. Non-constant-time field inversion is used to convert points from projective to affine in the precomputation table only, since it does not depend on the (secret) scalar.

7.2 Implementation timings

We start by describing our benchmarks for the Armv8 AArch64 implementation, written from scratch. We used the ODROID C4 microcontroller, as we wanted a smaller device that could be representative for the majority of Arm devices. It comes with a Quad-Core Cortex-A55, which is considered a mid-range processor. We employ `clang` from LLVM 13 with optimization level `-O3`.

Table 4. Benchmarks (in clock cycles) of the field arithmetic and elliptic curve point operations on an Arm Cortex-A55 2.0 GHz. The cost of reduction is included in the cost of field multiplication and squaring. Base field reduction is $\text{mod } zf(z)$. Op/m_b , Op/\tilde{m} denotes the cost of the operation relative to respectively base field and extension field multiplication.

Field op.	$\mathbb{F}_{2^{127}}$		$\mathbb{F}_{2^{254}}$	
	Cycles	Op/m_b	Cycles	Op/\tilde{m}
Multiplication	35	1.00	68	1.00
Reduction	16	0.46	15	0.22
Squaring	18	0.51	26	0.38
Inversion (ct.)	1 716	49.03	1 815	26.69
(non-ct.)	1 165	33.29	1 228	18.06

The benchmarks for our field implementation are presented in Table 4. Notice that non-constant time inversions that use lookup tables are roughly 33% faster.

Table 5. Constant-time variable base scalar multiplication benchmarks that are mostly performed on an Arm Quad-Core Cortex-A55 2.0 GHz. Memory is measured in terms of the number of elliptic curve points stored in the online precomputed table.

Implementation	Algorithm	Memory	Cycles
Lenngren [25] (Cortex-A55)	Curve25519	0	157,182
Longa [27] (Cortex-A55)	FourQ	8	191,184
Longa [27] (Cortex-A15)	FourQ	8	132,000
This work (Cortex-A55)	GLS254 1D $w = 5$	8	92,460
	GLS254 2D $w = 3$	4	86,525
	GLS254 2D $w = 4$	16	91,682

Table 5 presents our scalar multiplication timings in GLS254 and comparisons to related work. Compared to Intel, there are not a lot of efficient implementations specialized for Arm at the 128-bit security level. FourQ [10] is the closest competitor on Intel, and they also provide specialized implementations for Arm [27]. We benchmarked their Armv8 AArch64 implementation on our machine and included their Armv7 timings from [27] for the sake of fairness. A notable outlier is Lenngren’s implementation for Curve2559, which is a much closer competitor on Arm than any Curve25519 implementation on Intel.

As the first GLS254 implementation for Armv8, we are able to claim the constant-time scalar multiplication speed record by 34.5%. Contrary to the operation counts in Table 3, it is the 2D $w = 3$ algorithm that is the superior variant. This can be explained by the relatively high memory latency on our machine compared to higher-end models, which favors solutions that minimize the space used.

Table 6. Constant-time variable base scalar multiplication benchmarks for 64-bit Intel Core i7 4770 Haswell at 3.4GHz, and Core i7 7700 Kaby Lake at 3.6GHz, both with TurboBoost disabled. Memory is measured in terms of the number of elliptic curve points stored in the precomputed table.

Implementation	Algorithm	Memory	Cycles
Longa et al. [10] (Haswell)	FourQ	8	56,000
Longa et al. [10] (Kaby Lake)	FourQ	8	47,052
Oliveira et al. [30] (Haswell)	GLS254 1D $w = 5$	8	48,301
Oliveira et al. [31] (Skylake)	GLS254 1D $w = 5$	8	38,044
This work (Haswell)	GLS254 1D $w = 5$	8	45,966
	GLS254 2D $w = 3$	4	45,253
	GLS254 2D $w = 4$	16	47,184
This work (Kaby Lake)	GLS254 1D $w = 5$	8	36,480
	GLS254 2D $w = 3$	4	35,739
	GLS254 2D $w = 4$	16	38,076

For our Intel implementation, we extended the AVX-accelerated code from [34] with the new formulas and 2D variant. We do not report timings for field arithmetic due to space limitations, but they can be inferred directly from the timings in [30] with the speedups reported in [31]. We benchmarked our code in an older Core i7 4770 Haswell processor, and a Core i7 7700 Kaby Lake as the closest to the Skylake in [31]; both using `clang` from LLVM 13 and optimization level `-O3`. For 1D $w = 5$, we achieve a small speedup of 4.8% in Haswell and 4.1% for Skylake. The 2D $w = 3$ variant achieves a further speedup up to 2% using only half of the space. Surprisingly, the 2D $w = 4$ variant performs relatively poorly due to expensive conditional moves within the longer linear pass. The cumulative speedup over previous work is around 6% for both platforms. In comparison to FourQ, our timings are 24% faster and set a new speed record for constant-time scalar multiplication in Intel processors.

References

1. Ahmadi, O., Hankerson, D., Rodríguez-Henríquez, F.: Parallel formulations of scalar multiplication on koblitz curves. *J. Univers. Comput. Sci.* **14**(3), 481–504 (2008)
2. Azarderakhsh, R., Karabina, K.: A new double point multiplication algorithm and its application to binary elliptic curves with endomorphisms. *IEEE Trans. Computers* **63**(10), 2614–2619 (2014)
3. Bernstein, D.J.: Curve25519: New diffie-hellman speed records. In: *Public Key Cryptography. Lecture Notes in Computer Science*, vol. 3958, pp. 207–228. Springer (2006)
4. Bernstein, D.J., Lange, T.: SafeCurves: choosing safe curves for elliptic-curve cryptography. <https://safecurves.cr.yp.to/>
5. Bos, J.W., Costello, C., Hisil, H., Lauter, K.E.: High-performance scalar multiplication using 8-dimensional GLV/GLS decomposition. In: *CHES. LNCS*, vol. 8086, pp. 331–348. Springer (2013)
6. Bos, J.W., Costello, C., Longa, P., Naehrig, M.: Selecting elliptic curves for cryptography: an efficiency and security analysis. *J. Cryptogr. Eng.* **6**(4), 259–286 (2016)
7. Bos, J.W., Kleinjung, T., Niederhagen, R., Schwabe, P.: ECC2K-130 on cell cpus. In: *AFRICACRYPT. LNCS*, vol. 6055, pp. 225–242. Springer (2010)
8. Chi, J., Oliveira, T.: Attacking a binary GLS elliptic curve with magma. In: *LATINCRYPT. LNCS*, vol. 9230, pp. 308–326. Springer (2015)
9. Ciet, M., Lange, T., Sica, F., Quisquater, J.: Improved algorithms for efficient arithmetic on elliptic curves using fast endomorphisms. In: *EUROCRYPT. Lecture Notes in Computer Science*, vol. 2656, pp. 388–400. Springer (2003)
10. Costello, C., Longa, P.: Four \mathbb{Q} : four-dimensional decompositions on a \mathbb{Q} -curve over the Mersenne prime. *IACR Cryptol. ePrint Arch.* p. 565 (2015)
11. Faz-Hernández, A., Longa, P., Sánchez, A.H.: Efficient and secure algorithms for glv-based scalar multiplication and their implementation on GLV-GLS curves (extended version). *J. Cryptogr. Eng.* **5**(1), 31–52 (2015)
12. Galbraith, S.D., Gaudry, P.: Recent progress on the elliptic curve discrete logarithm problem. *Des. Codes Cryptogr.* **78**(1), 51–72 (2016)
13. Galbraith, S.D., Lin, X., Scott, M.: Endomorphisms for faster elliptic curve cryptography on a large class of curves. In: *EUROCRYPT. LNCS*, vol. 5479, pp. 518–535. Springer (2009)
14. Gallant, R.P., Lambert, R.J., Vanstone, S.A.: Faster point multiplication on elliptic curves with efficient endomorphisms. In: *CRYPTO. LNCS*, vol. 2139, pp. 190–200. Springer (2001)
15. Gaudry, P., Hess, F., Smart, N.P.: Constructive and destructive facets of weil descent on elliptic curves. *J. Cryptol.* **15**(1), 19–46 (2002)
16. Gouvêa, C.P.L., López-Hernández, J.C.: Implementing GCM on armv8. In: *CT-RSA. LNCS*, vol. 9048, pp. 167–180. Springer (2015)
17. Hankerson, D., Karabina, K., Menezes, A.: Analyzing the galbraith-lin-scott point multiplication method for elliptic curves over binary fields. *IEEE Trans. Computers* **58**(10), 1411–1420 (2009)
18. Hankerson, D., Vanstone, S., Menezes, A.: *Guide to Elliptic Curve Cryptography*. Springer-Verlag (2004)
19. Hess, F.: Generalising the ghs attack on the elliptic curve discrete logarithm problem. *LMS Journal of Computation and Mathematics* **7**, 167–192 (2004)

20. Itoh, T., Tsujii, S.: A fast algorithm for computing multiplicative inverses in $\text{gf}(2^m)$ using normal bases. *Inf. Comput.* **78**(3), 171–177 (1988)
21. Joye, M., Tunstall, M.: Exponent recoding and regular exponentiation algorithms. In: *AFRICACRYPT. LNCS*, vol. 5580, pp. 334–349. Springer (2009)
22. Kales, D., Rechberger, C., Schneider, T., Senker, M., Weinert, C.: Mobile private contact discovery at scale. In: *USENIX Security Symposium*. pp. 1447–1464. USENIX Association (2019)
23. Klaus Pommerening: Quadratic equations in finite fields of characteristic 2. <https://www.staff.uni-mainz.de/pommeren/MathMisc/QuGChar2.pdf> (2012)
24. Koblitz, A.H., Koblitz, N., Menezes, A.: Elliptic curve cryptography: The serpentine course of a paradigm shift. *Journal of Number theory* **131**(5), 781–814 (2011)
25. Lenngren, E.: AArch64 optimized implementation for X25519. <https://github.com/Emill/X25519-AArch64>
26. Lidl, R., Niederreiter, H.: *Finite fields*. Cambridge University Press (1997)
27. Longa, P.: FourQneon: Faster elliptic curve scalar multiplications on ARM processors. In: *SAC. LNCS*, vol. 10532, pp. 501–519. Springer (2016)
28. López-Hernández, J.C., Dahab, R.: New point compression algorithms for binary curves. In: *ITW*. pp. 126–130. IEEE (2006)
29. McLoughlin, M.B.: addchain: Cryptographic addition chain generation in go. Repository <https://github.com/mmcloughlin/addchain> (Oct 2021)
30. Oliveira, T., López-Hernández, J.C., Aranha, D.F., Rodríguez-Henríquez, F.: Two is the fastest prime: lambda coordinates for binary elliptic curves. *J. Cryptogr. Eng.* **4**(1), 3–17 (2014)
31. Oliveira, T., López-Hernández, J.C., Aranha, D.F., Rodríguez-Henríquez, F.: Improving the performance of the gls254. (2016). Presentation at CHES 2016 rump session (2016)
32. Oliveira, T., López-Hernández, J.C., Cervantes-Vázquez, D., Rodríguez-Henríquez, F.: Koblitz curves over quadratic fields. *J. Cryptol.* **32**(3), 867–894 (2019)
33. Renes, J., Costello, C., Batina, L.: Complete addition formulas for prime order elliptic curves. In: *EUROCRYPT (1)*. *Lecture Notes in Computer Science*, vol. 9665, pp. 403–428. Springer (2016)
34. Resende, A.C.D., Aranha, D.F.: Faster unbalanced private set intersection. In: *Financial Cryptography. Lecture Notes in Computer Science*, vol. 10957, pp. 203–221. Springer (2018)

Appendix

A Point compression for binary GLS curves

As a last construction, we present a new point compression algorithm for binary GLS curves for points in λ -affine form. The best known point compression algorithm for elliptic curves over \mathbb{F}_{2^n} is Algorithm 5 from [28]. It compresses an affine point (x, y) of $2n$ bits to n bits, and is the first to do so without needing an inversion. However, it requires $\text{Tr}(a) = 1$ and n to be odd. The latter condition is a problem for binary GLS curves, since they are defined over $\mathbb{F}_{2^{2m}}$ for an odd prime m . Our new algorithm adapts the techniques of [28] to this setting.

First we need some notation. Let $E = \mathbb{F}_{q^m}$ be the finite field extension of $K = \mathbb{F}_q$. Then the field trace $\text{Tr}_{E/K} : E \rightarrow K$ is defined as $\text{Tr}_{E/K}(c) = \sum_{i=0}^{m-1} c^{q^i}$ [26]. For our purposes, we define $q = 2^m$, $\text{Tr}' = \text{Tr}_{\mathbb{F}_{q^2}/\mathbb{F}_2}$ and $\text{Tr} = \text{Tr}_{\mathbb{F}_q/\mathbb{F}_2}$.

The point decomposition algorithm needs to solve a quadratic equation $\lambda^2 + \lambda = c$ in \mathbb{F}_{q^2} . The equation has a solution if and only if $\text{Tr}'(c) = 0$ [23, p. 54]. If a solution exists, it can be efficiently found using the technique from [17] that was generalized in [30]. Given a solution $\hat{\lambda}$, the other solution is $\hat{\lambda} + 1$.

Our algorithm works for any point $P = (x, \lambda)$ in the subgroup S of large prime order r . The compression algorithm computes $C_P = x + \text{lsb}(\lambda_0)u$ of m bits. Here $\text{lsb} : \mathbb{F}_q \rightarrow \mathbb{F}_2$ is the function that on input $d = d_0 + \dots + d_{q-1}z^{q-1} \in \mathbb{F}_q$ outputs d_0 . P can then be recovered from C_P as follows.

Algorithm 9: Decompression algorithm

Input : $C_P = x + \text{lsb}(\lambda_0)u$
Output: $P = (x, \lambda) \in S$

- 1 $t \leftarrow \text{Tr}(C_{P,1}) + 1$
- 2 $x \leftarrow c + tu$
- 3 Find solution λ' for $\lambda^2 + \lambda = b/x^2 + x^2 + a$
- 4 **if** $t = \text{lsb}(\lambda'_0)$ **then**
- 5 $\lambda \leftarrow \lambda'$
- 6 **else**
- 7 $\lambda \leftarrow \lambda' + 1$
- 8 **return** (x, λ)

Lemma 3. *Let $P = (x, \lambda) \in S$. Then Algorithm 9 recovers P from C_P .*

Proof. We are going to need some properties of the trace function. Firstly, $\text{Tr}_{E/K}$ is a linear transformation from E onto K , where E, K are viewed as vector spaces over K [23, p. 55]. The trace is also transitive, meaning for a finite extension F of E , $\text{Tr}_{F/K}(c) = \text{Tr}_{E/K}(\text{Tr}_{F/E}(c))$ [23, p. 56]. For binary elliptic curves it holds that for all $P = (x, \lambda)$ of odd order, $\text{Tr}'(x) = \text{Tr}'(a)$ [18, p. 130]. Finally, because m is odd, $\text{Tr}(d) = d$ for $d \in \mathbb{F}_2$.

From the transitivity of the trace function, we get that for $c = c_0 + c_1u \in \mathbb{F}_{q^2}$ that $\text{Tr}'(c) = \text{Tr}(c + c^q) = \text{Tr}((c_0 + c_1u) + (c_0 + c_1 + c_1u)) = \text{Tr}(c_1)$. As binary GLS curves have $\text{Tr}'(a) = 1$, it follows that $\text{Tr}(x) = 1$ for all $P = (x, \lambda) \in S$. Then $t = \text{Tr}(C_{P,1}) + 1 = (\text{Tr}(x_1) + \text{Tr}(\text{lsb}(\lambda_0))) + \text{Tr}(x_1) = \text{lsb}(\lambda_0)$.

Next, we correctly recover x as $C_P + tu$. The λ -affine curve equation is $\lambda^2x^2 + \lambda x^2 = x^4 + ax^2 + b$. Dividing both sides by x^2 , we get a quadratic equation in standard form, which can be solved using the technique from [17]. Note that $x \neq 0$ when P is represented in λ -affine coordinates. Given the candidate solutions $\hat{\lambda}$ and $\hat{\lambda} + 1$, we recover the correct one from t . \square

B The new formulas used for precomputation and exception-free execution

Proposition 1. *Let $P = (x, \lambda)$ be a point on $E'_{a,b}(\mathbb{F}_{q^2})$ with $3P \neq \mathcal{O}$. Then $3P$ can be computed in λ -projective coordinates as follows.*

$$\begin{aligned} T &= \lambda^2 + \lambda + a \\ A &= (x + T)^2 + T \\ X_{3P} &= x \cdot A^2 \\ Z_{3P} &= A \cdot (A + T) \\ \Lambda_{3P} &= T^2 \cdot T + Z_{3P}(\lambda + 1) \end{aligned}$$

Proposition 2. *Let $P = (x_P, \lambda_P)$ and $Q = (X_Q, \Lambda_Q, Z_Q)$ be points on $E'_{a,b}(\mathbb{F}_{q^2})$ with $P \neq \pm Q$. Then $P + Q$ and $P - Q$ can be simultaneously computed in λ -projective coordinates as follows.*

$$\begin{aligned} A &= \lambda_P \cdot Z_Q + \Lambda_Q \\ B &= (x_P \cdot Z_Q + X_Q)^2 \\ C &= (x_P \cdot Z_Q) \cdot X_Q \\ X_{P+Q} &= A^2 \cdot C \\ Z_{P+Q} &= A \cdot B \cdot Z_Q \\ \Lambda_{P+Q} &= (A \cdot X_Q + B)^2 + Z_{P+Q} \cdot (\lambda_P + 1) \\ X_{P-Q} &= X_{P+Q} + C \cdot Z_Q^2 \\ Z_{P-Q} &= Z_{P+Q} + B \cdot Z_Q^2 \\ \Lambda_{P-Q} &= \Lambda_{P+Q} + (X_Q \cdot Z_Q)^2 + (B \cdot Z_Q^2) \cdot (\lambda_P + 1) \end{aligned}$$

Proposition 3. *Let $P = (X, \Lambda, Z)$ on $E'_{a,b}(\mathbb{F}_{q^2})$. For $c \in \mathbb{F}_{q^2}$, let c_0, c_1 denote its coefficients in \mathbb{F}_q such that $c = c_0 + c_1u$. Then $\psi(P)$ can be computed in λ -projective coordinates as follows.*

$$\begin{aligned} X_{\psi(P)} &= X + X_1 \\ \Lambda_{\psi(P)} &= \Lambda + \Lambda_1 + Z_1 + Z_0u \\ Z_{\psi(P)} &= Z + Z_1 \end{aligned}$$

Proposition 4. *Let $P = (x, \lambda)$ on $E'_{a,b}(\mathbb{F}_{q^2})$. For $c \in \mathbb{F}_{q^2}$, let c_0, c_1 denote its coefficients in \mathbb{F}_q such that $c = c_0 + c_1u$. Then $P + \psi(P)$ can be computed in λ -projective coordinates as follows.*

$$\begin{aligned}
A &= (x_0 \cdot \lambda_1 + x_1) + (x_1 \cdot \lambda_1 + x_0 + x_1)u \\
B &= x_1 \cdot \lambda_1 + x_1u \\
X_{P+\psi(P)} &= A \cdot (A + B) \\
Z_{P+\psi(P)} &= (x_1^2 \cdot \lambda_1) + x_1^2u \\
\Lambda_{P+\psi(P)} &= (A + B + x_1^2)^2 + Z_{P+\psi(P)} \cdot (\lambda_P + 1)
\end{aligned}$$

Finally, we also include how to compute $2Q + P$ and $2Q + P_1 + P_2$ for respectively the 1D and 2D variant in the last iteration of the scalar multiplication loop, using complete formulas. Our approach was to compute them non-atomically from left to right. We do a normal doubling which is exception-free and use Algorithm 10 for the complete mixed additions. A complete mixed addition costs $11\tilde{m} + 5\tilde{s}$, which is slightly more than a full addition. To make Algorithm 10 run in constant time, CMOV instructions are used for the conditional assignments.

Algorithm 10: Complete mixed addition

Input : $P = (x_P, \lambda_P)$ in λ -affine coordinates, $Q = (X_Q, \Lambda_Q, Z_Q)$ in λ -projective coordinates.

Output: $P + Q$ in λ -projective coordinates.

- 1 $R \leftarrow P + Q$ using the incomplete formula.
- 2 $R_D \leftarrow 2P$
- 3 $X_P \leftarrow x_P \cdot Z_Q$
- 4 $\Lambda_P \leftarrow \lambda_P \cdot Z_Q$
- 5 **if** $Z_Q = 0$ **then**
- 6 $R \leftarrow P$
- 7 **if** $X_P = X_Q$ **and** $\Lambda_P = \Lambda_Q$ **then**
- 8 $R \leftarrow R_D$
- 9 **if** $X_P = X_Q$ **and** $\Lambda_P = \Lambda_Q + Z_Q$ **then**
- 10 $R \leftarrow \mathcal{O}$
- 11 **return** R
