

Efficient Algorithms for Large Prime Characteristic Fields and Their Application to Bilinear Pairings

Patrick Longa

Microsoft Research, USA
plonga@microsoft.com

Abstract. We propose a novel approach that generalizes interleaved modular multiplication algorithms for the computation of sums of products over large prime fields. This operation has widespread use and is at the core of many cryptographic applications. The method reformulates the widely used lazy reduction technique, crucially avoiding the need for storage and computation of “double-precision” operations. Moreover, it can be easily adapted to the different methods that exist to compute modular multiplication, producing algorithms that are significantly more efficient and memory-friendly. We showcase the performance of the proposed approach in the computation of multiplication over an extension field \mathbb{F}_{p^k} , and demonstrate its impact with record-breaking implementations of bilinear pairings. Specifically, we accomplish a full optimal ate pairing computation over the popular BLS12-381 curve, designed for the 128-bit security level, in under half a millisecond on a 3.2GHz Intel Coffee Lake processor, which is about $1.40\times$ faster than the state-of-the-art. Similarly, we perform the same computation over the BLS24-509 curve, targeting the 192-bit security level, in ~ 2.6 milliseconds, achieving a speedup of more than $1.30\times$. We also report a significant impact on other applications, including protocols based on supersingular isogenies.

Keywords: Sum of products · prime fields · extension fields · bilinear pairings · BLS12-381 · supersingular isogenies · efficient computation.

1 Introduction

Take two sets of integers $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ all defined over a certain finite field \mathbb{F}_p with large prime characteristic p . The sum of their products, namely, the computation $c = \sum_{i=0}^{t-1} \pm a_i \cdot b_i \bmod p$ is a fundamental operation that is at the core of various types of arithmetic over \mathbb{F}_p , from matrix multiplication and multiplication over polynomial rings to elliptic curve arithmetic. Of special interest is that this operation has wide use in the form of multiplication over extension fields \mathbb{F}_{p^k} , which is at the heart of several cryptographic schemes such as elliptic curves defined over extension fields [37], bilinear pairings on ordinary curves [71] and supersingular isogeny-based schemes [55].

In this work, we propose a new approach that computes a sum of products over \mathbb{F}_p by interleaving intermediate products with the modular reduction step, in similar fashion to classical interleaved modular multiplication algorithms (§2.1). Crucially, it departs from algorithms using the well-known lazy reduction technique [74, 78], eliminating the excessive growth of intermediate values and the need of computing “double-precision” arithmetic. The method can be easily adapted to existing algorithmic variants that fit different application profiles, for software and hardware platforms. We show that some of these variants are especially efficient for software implementation, exhibiting strong synergy

with computer architectures that leverage the simplicity of schoolbook-like multiplication. The resulting outcome is streamlined implementations that operate with significantly reduced memory usage.

The target. To showcase the potential of the new method, we apply it in the context of bilinear pairings, which are at the core of a myriad of elegant cryptographic schemes such as identity-based cryptosystems [20, 24, 29, 30] and non-interactive zero-knowledge proof systems [52, 53]. More recently, pairings have attracted great interest because of their novel use in blockchain technology like Zcash [6] and Ethereum 2.0 [27]. Unfortunately, recent advances in the computation of discrete logarithms over extension fields, as those used in pairing-based cryptography [10, 58], have forced an increase in parameter sizes [9, 68]. This, for example, motivated the design of BLS12-381, a new pairing-friendly curve that is conjectured to provide close to 128 bits of security and is widely used in blockchain protocols [23] (see [73] for details on a standardization effort in the IETF CFRG). The increased parameter sizes, coupled with the inherently computationally expensive nature of pairing operations, underscore the need of devising innovative methods that improve performance.

We note that the proposed approach extends beyond the efficient implementation of multiplication over extension fields (e.g., it could also be used for elliptic curve arithmetic operations with the form $AB+CD$ over a finite field of large prime characteristic). However, as we anticipate more substantial gains for arithmetic over extension fields, we use this setting within the context of pairings for illustrative purposes. Other attractive applications include post-quantum supersingular isogeny-based signature schemes such as SQISign [45] and supersingular isogeny-based zero-knowledge proof schemes such as the protocol by Basso et al. [17]¹; see §5.1 for more details.

Despite its broader generality, the method is analyzed in the context of Montgomery multiplication [69], as it is the most prevalent approach for modular multiplication.

Computing sums of products and the case of multiplication over \mathbb{F}_{p^k} . There are two main approaches to optimize multiplication over \mathbb{F}_{p^k} . On one hand, algebraic transformations are used to reduce the required number of underlying field multiplications. A well-known example of this case is Karatsuba multiplication [57].

The second approach consists in minimizing the number of modular reductions using the so-called *lazy reduction* technique. Lazy reduction, which goes back to at least [78], is an extensively used optimization that has been applied in a wide variety of scenarios [4, 31, 62, 65, 74]. The basic principle is very simple: products are computed and left unreduced. A modular reduction is only applied at the very end of the computation, right after the summation of the intermediate, “double-precision” values. This elimination of reductions is highly effective, especially for primes for which the reduction routine is roughly as expensive as the integer multiplication part. If we assume the use of Montgomery multiplication for computing the summation of t n -digit products (§2.2), the cost is reduced from $t(2n^2 + n)$ multiplies to only $n^2(t + 1) + n$.

In the context of cryptographic pairings, Scott [74] was the first to apply lazy reduction to Karatsuba multiplication in the computation of multiplication over \mathbb{F}_{p^2} . Later, the technique was extended to the full tower and curve arithmetic by Aranha et al. [3] (see also [66]). Given an extension field $\mathbb{F}_{p^k} = \mathbb{F}_p[x]/(x^k - \omega)$ with $\omega \in \mathbb{F}_p$, a multiplication in \mathbb{F}_{p^k} exploiting lazy reduction can be performed with only k reductions modulo p , in contrast to the k^2 reductions that would be required with a conventional multiplication.

¹Other very attractive applications where the method would have had a significant impact include SIDH [55] and SIKE [5]; see Appendix B. Unfortunately, these protocols were very recently proven insecure in a series of papers starting with [28].

Many works have exploited the technique (especially in the context of multiplication over \mathbb{F}_{p^2}) without changing the basic approach [3, 7, 19, 21, 22, 38, 76]. In fact, despite the availability of efficient interleaved modular multiplication algorithms [42], the combination of lazy reduction with these faster algorithms has been long seen as an impossibility. As consequence, lazy reduction has been strictly limited to use non-interleaved modular multiplications.

The main drawback of this traditional approach using lazy reduction is that it forces the storage and computation with intermediate results of double precision. In addition to the increased pressure on memory usage, the optimization of double-precision arithmetic may require specialized routines that significantly increase the implementation complexity [2, 3]. In contrast, the proposed method limits the growth of intermediate results and gets rid of double-precision arithmetic. The simplest variants use schoolbook internally, which eliminates further the additional storage demanded by Karatsuba. For small and medium-sized primes, the significant reduction in memory friction pays off in terms of speed, even in the cases in which the use of multiplications is higher, as we show in §4.

Open-source software. We have implemented and integrated our algorithms to RELIC [2], a cryptographic library containing state-of-the-art implementations of pairings. Our implementations cover two pairing-friendly curves: BLS12-381 and BLS24-509, which are intended to match the 128- and 192-bit security levels, respectively. Both are instantiated with an optimal ate pairing. The software is available as open-source at:

https://github.com/primefieldsfp/faster_Fpx.

The software stack also includes the code that we wrote to evaluate the speed performance and memory usage of the proposed method using the SIKE primes, as described in Appendix B.

Outline. We start by giving some preliminary background on algorithmic aspects of Montgomery multiplication and extension field arithmetic in §2. Then, we describe the details of our method in §3, together with an exhaustive classification of its different algorithmic variants. This section also includes a preliminary cost analysis. In §4, we present our case study targeting pairings, describe various implementation choices for different extension fields, and report efficient implementations of the BLS12-381 and BLS24-509 pairing-friendly curves using the RELIC library. Finally, we discuss potential future developments and the impact of this work, including the case of supersingular isogeny-based protocols, in §5.

2 Preliminaries

2.1 Montgomery multiplication

A well-known and widely-used method to implement modular multiplication is due to Montgomery [69]. This method introduces a significant speedup in the computation of modular reduction by replacing expensive long divisions by simple divisions with powers of two.

To carry out a Montgomery multiplication, field elements are represented in the so-called Montgomery domain. Let $R = 2^N$ and $p' = -p^{-1} \bmod R$, where $N = n \cdot w$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$ and w is the computer wordsize. For two field elements a and b , their Montgomery representation is given by $\tilde{a} = aR \bmod p$ and $\tilde{b} = bR \bmod p$, respectively. If it holds that $\tilde{a}\tilde{b} < pR$, the Montgomery residue $c = \tilde{a}\tilde{b}R^{-1} \bmod p = abR \bmod p$ is then computed as

$$c = (\tilde{a}\tilde{b} + (\tilde{a}\tilde{b}p' \bmod 2^N) \cdot p) / 2^N. \quad (1)$$

The final result $a \cdot b$ can then be easily obtained by dividing by the value R . If one assumes that conversions to/from the Montgomery domain are amortized by executing a long series of modular multiplications, the cost of a Montgomery reduction (i.e., without the integer multiplication part) is given by $(n^2 + n)$ w -bit multiplications.

Interleaved and non-interleaved modular multiplication There are two general approaches for implementing modular multiplication, which are determined by how the integer multiplication and reduction parts are “glued” together. In an *interleaved* or *non-separated* modular multiplication, both pieces are computed in an intertwined fashion, while the *non-interleaved* or *separated* version computes the full integer multiplication first and then proceeds to carry out the reduction part separately. In applications dealing with sums of products over large prime fields (e.g., for multiplying over extension fields like in pairings or supersingular isogeny-based protocols), the use of non-interleaved modular multiplication has become the *de facto* approach, as it enables a straightforward application of lazy reduction and other advanced implementation techniques [3, 4, 19, 21, 22, 38, 62, 74, 76, 78].

Radix- r Montgomery multiplication. Straight implementations of Eq. (1) demand the use of a high number of registers since the full inputs are processed in a single pass using the full modulus. Let r be a certain radix, typically a power of two, in which operands and the modulus are represented (e.g., an operand a is represented as $(a_{t-1}, \dots, a_1, a_0)_r$). A more implementation-friendly approach proposed by Dussé and Kaliski Jr. [42] processes the computation one digit at a time reducing with r at each iteration, in what is called the radix- r Montgomery reduction. An *interleaved* computation of a radix- r Montgomery multiplication of two Montgomery elements \tilde{a} and \tilde{b} then proceeds by fixing $p' = -p^{-1} \bmod r$ (assuming that the modulus is a prime p of bitlength l), initializing c to 0, and computing $t = \lceil l / \log_2 r \rceil$ iterations doing

$$c = (c + \tilde{a}_i \tilde{b} + ((c + \tilde{a}_i \tilde{b})p' \bmod r) \cdot p) / r, \quad (2)$$

for $i = 0, \dots, t - 1$.

In this work, we adopt a generalization of the original radix- r Montgomery multiplication by setting $r = 2^{Bw}$, where $B \in \mathbb{Z}$ and $0 < B \leq n$ and, as before, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$ and w is the computer wordsize². In this case, Eq. (2) runs for $\lceil n/B \rceil$ iterations. This generalization lifts the restriction that the bitlength of the radix r should match the computer wordsize w of a given platform, as originally assumed in [42]. Note that the original radix- r Montgomery multiplication corresponds to the case $B = 1$.

As we will show in §3, the flexibility introduced by B in the definition of the radix allows for a comprehensive generalization that captures many implementation variants of Montgomery multiplication exploiting either the “operand-scanning” form (a.k.a. school-book method), the “product-scanning” form (a.k.a. Comba method [33]), the Karatsuba technique [57], and their different combinations. To the best of our knowledge, several of the arising variants are novel.

2.2 Sums of products over large prime fields

Let $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ be two sets of elements all belonging to a certain field \mathbb{F}_p of large prime characteristic p . We define a “sum of products” as a computation with the form

$$c = \sum_{i=0}^{t-1} \pm a_i \cdot b_i \bmod p$$

²This generalization is similar to the description by Bos and Friedberger [21, §3.2], but without limiting to a special-form prime.

This operation can be found at the core of many cryptologic computations, the most notable of which is perhaps multiplication over extension fields with large prime characteristic. Hence, we use this operation to illustrate the impact of the proposed algorithms (see §3).

Let's recall the simplest scenario for a multiplication of two elements $a, b \in \mathbb{F}_{p^k}$ modulo an irreducible polynomial with the form $f = x^k - \omega$, where ω is a primitive element in \mathbb{F}_p^* and $k|(p-1)$. The polynomial multiplication is then given by

$$\begin{aligned} c(x) = a(x)b(x) &= \left(\sum_{i=0}^{k-1} a_i x^i \right) \left(\sum_{i=0}^{k-1} b_i x^i \right) \\ &\equiv c_{k-1} x^{k-1} + \sum_{i=0}^{k-2} (c_i + \omega c_{i+k}) x^i \pmod{f(x)}, \end{aligned}$$

where

$$c_j = \sum_{i=0}^j a_i b_{j-i} \pmod{p}.$$

It is a widespread practice to optimize the implementation for computing each c_j using an ‘‘accumulation and reduction’’ strategy, most commonly known as lazy reduction. This technique effectively reduces the number of modular reductions to only one (or k , for the full polynomial multiplication). Note that, as in most practical scenarios, we assume that ω has small coefficients that make a multiplication by it relatively cheap.

As mentioned before, this use of lazy reduction has some drawbacks, the most critical of which being the need of extra storage and computing with intermediate results of double precision. As we show in §4 (see also Appendix B), this issue makes implementations slower and less memory-friendly for small and medium-sized primes.

The rise of fast multiplication over \mathbb{F}_{p^2} . The most basic ‘‘sum of products’’ operation underlying several cryptographic schemes is multiplication over a quadratic extension field \mathbb{F}_{p^2} . Modern examples of these schemes include bilinear pairings on ordinary elliptic curves over prime fields and supersingular isogeny-based protocols.

For illustrative purposes, let's use the common construction fixing $\mathbb{F}_{p^2} = \mathbb{F}_p(i)$ for $i^2 - \beta = 0$, where β is a small integer in absolute value. Two main approaches are known to realize the multiplication in this case.

Take two elements $a = (a_0 + a_1 i)$ and $b = (b_0 + b_1 i) \in \mathbb{F}_{p^2}$. The first method to compute the multiplication $a \cdot b$ in \mathbb{F}_{p^2} is the straightforward schoolbook method which computes it as

$$a \cdot b = (a_0 b_0 + a_1 b_1 \beta) + (a_0 b_1 + a_1 b_0) i$$

The second approach is Karatsuba multiplication, which computes the same operation as

$$a \cdot b = (a_0 b_0 + a_1 b_1 \beta) + ((a_0 + a_1)(b_0 + b_1) - a_0 b_0 - a_1 b_1) i$$

If we assume that $\beta = -1$ as in most efficient instantiations, the operation count for multiplication using schoolbook is of four modular multiplications and two modular additions/subtractions, while the one for Karatsuba is of three modular multiplications and five modular additions/subtractions.

Efficient implementation of the arithmetic over \mathbb{F}_{p^2} attracted lots of interest from the cryptographic community around the mid-2000's, contributing to a remarkable effort aimed

at reducing the high cost of computing cryptographic pairings. In 2007, Scott [74] was the first to apply lazy reduction to Karatsuba multiplication in the context of pairings, changing the cost of one \mathbb{F}_{p^2} multiplication to three integer multiplications, two modular reductions, two additions and three double-precision additions/subtractions. This approach was later perfected by Beuchat et al. [19] and Aranha et al. [3] with the use of some aggressive optimization techniques such as the avoidance of modular corrections and carry-handling elimination, all in the context of the computation of optimal ate pairings [77] over a 254-bit Barreto–Naehrig (BN) curve [15].

The algorithm for multiplication in \mathbb{F}_{p^2} combining all these optimizations for the optimal case with $\beta = -1$ is shown in Algorithm 1. Integer operations without modular correction or reduction are represented as \times , $+$ or $-$. The only operation that requires a modular correction is the subtraction in line 8 that is represented as \ominus . Double-precision operands are represented in uppercase, while single-precision operands are in lowercase.

Algorithm 1 Multiplication in \mathbb{F}_{p^2} using Karatsuba and lazy reduction

Input: $a = (a_0 + a_1i)$ and $b = (b_0 + b_1i) \in \mathbb{F}_{p^2}$

Output: $c = a \cdot b = (c_0 + c_1i) \in \mathbb{F}_{p^2}$

1: $T_0 \leftarrow a_0 \times b_0$ 2: $T_1 \leftarrow a_1 \times b_1$ 3: $t_0 \leftarrow a_0 + a_1$ 4: $t_1 \leftarrow b_0 + b_1$ 5: $T_2 \leftarrow t_0 \times t_1$ 6: $T_3 \leftarrow T_0 + T_1$	7: $T_2 \leftarrow T_2 - T_3$ 8: $T_0 \leftarrow T_0 \ominus T_1$ 9: $c_0 \leftarrow T_2 \bmod p$ 10: $c_1 \leftarrow T_0 \bmod p$ 11: return $C = (c_0 + c_1i)$
---	---

The case of higher-degree field extensions. For pairings, high-degree extension field arithmetic represents the main building block and, therefore, its efficient implementation becomes crucial. To this end, it has been recommended to implement it as a tower of extensions built with suitable irreducible polynomials [60], following a similar development for optimal extension fields [8]. For example, Pereira et al. [50] recommended to set $p^e \equiv 1 \pmod{6}$ and represent $\mathbb{F}_{p^{6e}}$ as a tower extension of \mathbb{F}_{p^e} in one of three ways:

- $\mathbb{F}_{p^{6e}} = \mathbb{F}_{p^e}[u]/(u^6 - \xi)$,
- $\mathbb{F}_{p^{6e}} = \mathbb{F}_{p^{2e}}[v]/(v^3 - \xi)$ with $\mathbb{F}_{p^{2e}} = \mathbb{F}_{p^e}[s]/(s^2 - \xi)$,
- $\mathbb{F}_{p^{6e}} = \mathbb{F}_{p^{3e}}[w]/(w^2 - \xi)$ with $\mathbb{F}_{p^{3e}} = \mathbb{F}_{p^e}[t]/(t^3 - \xi)$,

where $\xi \in \mathbb{F}_{p^e}$ is a non-square and a non-cube.

The use of a tower field allows to write modularized implementations, in which each layer can be easily optimized using algebraic transformations like Karatsuba to reduce the number of modular multiplications.

In 2011, Aranha et al. [3] pushed the performance limits further by extending the use of lazy reduction to the full tower scheme, in order to minimize the use of modular reductions. Concretely, they showed that this optimization, when applied above the \mathbb{F}_{p^2} arithmetic up to the $\mathbb{F}_{p^{12}}$ layer, results in an 11%-17% speedup on a variety of x64 processors. Nevertheless, this extended lazy reduction technique comes at a price. It requires additional specialized routines to perform the double-precision arithmetic, which increase the complexity and memory footprint of the implementation.

In §4, we discuss how one can improve performance and memory use for prime sizes of practical relevance with a new approach that we present next.

3 The Proposed Method

We describe the proposed method in the context of Montgomery multiplication, which is arguably one of the most relevant scenarios. See §5 for a discussion on other potential applications.

Let $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ be two sets of integers with $a_i, b_i \in [0, p)$ for $i = 0, \dots, (t-1)$ and $0 \leq \sum_{i=0}^{t-1} a_i b_i < pR$, where $R = 2^{nw}$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$, and w is the computer wordsize. From now on, we make the assumption that inputs a_i and b_i are already in the Montgomery domain.

From a general perspective, the new approach essentially consists in performing a *merged* computation of the operation $c = \sum_{i=0}^{t-1} \pm a_i \cdot b_i \bmod p$ using *interleaved* radix- r Montgomery multiplication³, that is, initializing $c = 0$ and executing $\lceil l/\log_2 r \rceil = \lceil n/B \rceil$ iterations doing

$$c = \left(c + \sum_{i=0}^{t-1} a_{i,j} b_i + \left(\left(c + \sum_{i=0}^{t-1} a_{i,j} b_i \right) p' \bmod r \right) \cdot p \right) / r, \quad (3)$$

for $j = 0, \dots, \lceil n/B \rceil - 1$, where $p' = -p^{-1} \bmod r$, and each integer a_i is represented in radix- r representation as $(a_{i,\lceil n/B \rceil - 1}, \dots, a_{i,1}, a_{i,0})_r$. As stated in §2.1, the radix r is adopted in the generalized form $r = 2^{Bw}$, where $B \in \mathbb{Z}$ and $0 < B \leq n$. In the following, we call each digit in this radix representation a “B-digit”.

We remark that Eq. (3) is presented in a general form for simplicity purposes. Next, we provide a more detailed description that covers the wide diversity of variants that can be derived from the approach.

At a high-level, we can classify the different variants by the method that is used to implement the top layer in the computation of the multiplications in Eq. (3). Thus, we can distinguish operand-scanning (or schoolbook), product-scanning (or Comba), and Karatsuba variants. In the remainder, we mostly focus on the first case which brings very fast computations to the software platform class that we target in this work. We comment that product-scanning and Karatsuba variants, such as those described in Appendix A, might be useful in other scenarios, e.g., for hardware implementations (see discussion in §5).

Remark 1. The result of a Montgomery reduction is upper bounded by $2p$ when its input is in the range $[0, pR)$. Hence, a conditional subtraction is needed to bring the result to $[0, p)$. However, this operation can be avoided if we perform arithmetic over a redundant representation (e.g., over \mathbb{Z}_{2p}). For example, if operands are kept in the range $[0, 2p)$ such that the result of a multiplication is guaranteed to be $c = a \cdot b < 4p^2 \leq pR$ (i.e., it should hold that $R \geq 4p$), then the result of the Montgomery reduction is still going to be bounded by $2p$ but we will no longer require the modular correction. A simple correction is going to be required at the very end of the computations to bring the final result to the canonical range $[0, p)$. In the following, all the algorithms assume the use of this redundant representation to avoid the final conditional subtraction.

3.1 Operand-scanning method

For this method, the computation flow at the top layer follows the operand-scanning or schoolbook form. That is, for each multiplication, a B -digit from the radix- r representation of a given multiplier is first multiplied with the full value of the multiplicand before proceeding to the next computation. For the remainder, we refer to this operation as *B-digit* \times *row* multiplication.

³As explained before, the non-interleaved or separated case is used with the standard lazy reduction technique.

We distinguish two main approaches, depending on whether the inner multiplications $a_{i,j}b_i$ from Eq. (3) are interleaved with the multiplications with the prime p or not. We adopt the naming convention from [61] and call the approaches *finely integrated* if we do the former case (i.e., with interleaved inner multiplications), and *coarsely integrated*, otherwise.

Coarsely integrated variants. The merged “sum of products” algorithm using a coarsely integrated strategy is displayed in Algorithm 2. The construction of the algorithm easily follows from Eq. (3) when $n \bmod B = 0$. We still need to manage the cases in which $n \bmod B \neq 0$ (i.e., the digit-size of the most significant B -digit is strictly smaller than that of a B -digit). This is done in lines 6–9, where the computations are adjusted to the right digit size.

As can be seen, the B -digit \times row multiplications corresponding to $a_{i,j}b_i$ (line 3) are interleaved with those for the multiplications with p' and p (lines 4 and 5) at each iteration of the for-loop.

There are multiple ways in which the inner multiply-and-accumulate operations $\sum a_{i,j} \cdot b_i$ and $u + q \cdot p$ can be realized. We classify these variants according to the chosen value B as follows:

- Case with $B = 1$: one is setting $r = 2^w$ and all the inner computations essentially become straight *digit* \times *row* multiplications. This is the analogous version of “Improvement 2” from [42], called coarsely integrated operand scanning (CIOS) in [61].
- Case with $B > 1, B \neq n$: the inner computations work on “blocks” of digits and, hence, each B -digit \times B -digit multiplication can be implemented in either schoolbook, Comba or Karatsuba style (or any combination of these in a *multi-level* fashion for sufficiently large primes).
- Case with $B = n$: this is essentially the original lazy reduction technique.

Algorithm 2 Merged sums of products using radix- r interleaved Montgomery multiplication in coarsely integrated form.

Input: integers $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ s.t. $a_i, b_i \in [0, 2p)$ for $i = 0, \dots, (t-1)$ and $0 \leq \sum_{i=0}^{t-1} a_i b_i < pR$, where $R = 2^{nw}$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$, and w is the computer wordsize; the radix $r = 2^{Bw}$ s.t. $B \in \mathbb{Z}$ and $0 < B \leq n$, and the Montgomery constant $p' = -p^{-1} \bmod r$. Integers are represented in radix r , e.g., $a_i = (a_{i, \lceil n/B \rceil - 1}, \dots, a_{i,1}, a_{i,0})_r$.

Output: the Montgomery residue $c = \sum_{i=0}^{t-1} a_i b_i \cdot R^{-1} \bmod p$ s.t. $c \in [0, 2p)$.

```

1:  $u \leftarrow 0$ 
2: for  $j = 0$  to  $\lfloor n/B \rfloor - 1$  do
3:    $u \leftarrow u + \sum_{i=0}^{t-1} a_{i,j} \cdot b_i$ 
4:    $q \leftarrow u \cdot p' \bmod 2^{Bw}$ 
5:    $u \leftarrow (u + q \cdot p) / 2^{Bw}$ 
6: if  $n \bmod B \neq 0$  then
7:    $u \leftarrow u + \sum_{i=0}^{t-1} a_{i, \lceil n/B \rceil - 1} \cdot b_i$ 
8:    $q \leftarrow u \cdot p' \bmod 2^{(n \bmod B)w}$ 
9:    $u \leftarrow (u + q \cdot p) / 2^{(n \bmod B)w}$ 
10: return  $c \leftarrow u$ 

```

Finely integrated variants. The merged “sum of products” algorithm using a finely integrated strategy is displayed in Algorithm 3. In this case, note that multiplications are performed B -digit by B -digit, interleaving those corresponding to $a_{i,k}b_{i,j}$ (lines 3 and 7) with those for the multiplications with p_k (lines 5 and 8). Note that we manage the case with $n \bmod B \neq 0$ as described before.

Similar to the coarsely integrated form, there are multiple ways in which the inner multiply-and-accumulate operations can be realized. Again, we classify these variants according to the value B as follows:

- Case with $B = 1$: one is setting $r = 2^w$ and all the inner computations become simple $digit \times digit$ multiplications, where those corresponding to input operands are interleaved with those with the prime. This is the analogous version of the finely integrated operand scanning (FIOS) method from [61].
- Case with $B > 1, B \neq n$: the inner computations work on “blocks” of digits and, hence, each B -digit \times B -digit multiplication can be implemented in either schoolbook, Comba or Karatsuba style (or any combination of these in a *multi-level* fashion for sufficiently large primes).
- Case with $B = n$: this is essentially the original lazy reduction technique.

Algorithm 3 Merged sums of products using radix- r interleaved Montgomery multiplication in finely integrated form.

Input: integers $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ s.t. $a_i, b_i \in [0, 2p)$ for $i = 0, \dots, (t-1)$ and $0 \leq \sum_{i=0}^{t-1} a_i b_i < pR$, where $R = 2^{nw}$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$, and w is the computer wordsize; the radix $r = 2^{Bw}$ s.t. $B \in \mathbb{Z}$ and $0 < B \leq n$, and the Montgomery constant $p' = -p^{-1} \bmod r$. Integers are represented in radix r , e.g., $a_i = (a_{i, \lceil n/B \rceil - 1}, \dots, a_{i,1}, a_{i,0})_r$.
Output: the Montgomery residue $c = \sum_{i=0}^{t-1} a_i b_i \cdot R^{-1} \bmod p$ s.t. $c \in [0, 2p)$.

```

1:  $u \leftarrow 0$ 
2: for  $j = 0$  to  $\lfloor n/B \rfloor - 1$  do
3:    $u \leftarrow u + \sum_{i=0}^{t-1} a_{i,0} \cdot b_{i,j}$ 
4:    $q \leftarrow u \cdot p' \bmod 2^{Bw}$ 
5:    $u \leftarrow (u + q \cdot p_0) / 2^{Bw}$ 
6:   for  $k = 1$  to  $\lceil n/B \rceil - 1$  do
7:      $u \leftarrow u + 2^{(k-1)Bw} \cdot \sum_{i=0}^{t-1} a_{i,k} \cdot b_{i,j}$ 
8:      $u \leftarrow u + 2^{(k-1)Bw} \cdot q \cdot p_k$ 
9: if  $n \bmod B \neq 0$  then
10:   $u \leftarrow u + \sum_{i=0}^{t-1} a_{i,0} \cdot b_{i, \lceil n/B \rceil - 1}$ 
11:   $q \leftarrow u \cdot p' \bmod 2^{(n \bmod B)w}$ 
12:   $u \leftarrow (u + q \cdot p_0) / 2^{(n \bmod B)w}$ 
13:  for  $k = 1$  to  $\lceil n/B \rceil - 1$  do
14:     $u \leftarrow u + 2^{(kB - n \bmod B)w} \cdot \sum_{i=0}^{t-1} a_{i,k} \cdot b_{i, \lceil n/B \rceil - 1}$ 
15:     $u \leftarrow u + 2^{(kB - n \bmod B)w} \cdot q \cdot p_k$ 
16: return  $c \leftarrow u$ 

```

Selecting a variant. Picking a specific variant depends on both the modulus size (see the next subsection) and the targeted platform. Generally speaking, the coarsely integrated variant (Algorithm 2) should be the preferred option in most software platforms in which schoolbook works well and the availability of general purpose registers (GPRs) is sufficient to support a full B -digit \times row multiplication with minimal interaction with the memory.

On the other hand, the inner for-loop of the finely integrated variant (Algorithm 3) consists of a bunch of multiplications that are independent from each other and, hence, can be executed in parallel on, e.g., an FPGA. See §5 for an extended discussion on other uses for the algorithms.

Regarding the selection of the value B , setting $B > 1$ might make it easier to alleviate memory use for relatively large primes, especially in the case of Algorithm 3. For small and medium size primes⁴, it appears that setting $B = 1$ hits the right balance between the size of intermediate results and the number of GPRs available on x64 platforms.

Finally, we mentioned before that for the internal multiplications it is possible to use either schoolbook, Comba, Karatsuba, or any combination of these in a multi-layer implementation. To be efficient, Karatsuba would require a B -digit consisting of a sufficiently large number of limbs. And between schoolbook and Comba, the former is typically preferable when a given platform supports efficient carry-saving instructions such as `mulx` or versatile multiply-and-add (`MULADD`) instructions (see §5).

Primes of special form. The algorithms described in this section have been presented generically assuming primes with no special form, as typically found in, e.g., pairing computations (§4). However, they can be easily adapted to settings using primes with some special shape. In particular, a relevant scenario is when $p+1 \equiv 0 \pmod{2^{zw}}$ for some integer $z \geq 1$. As an example of this case, see the adaptation of Algorithm 2 to the case of the SIKE primes in Appendix B. It is straightforward to apply a similar optimization to primes such as the SQISign prime p_{3923} [46], as discussed in §5.1.

We note that, in part, the proposed method gains in efficiency by eliminating modular reductions. Therefore, it is naturally expected that primes that support a very fast reduction (e.g., Mersenne or pseudo-Mersenne primes) do not gain a significant advantage with the approach.

3.2 Cost analysis

Except for the variants that could use Karatsuba at the lower levels of their computations (which would only be the case for relatively large primes), the complexity of the proposed algorithms is quadratic in terms of multiplication instructions. For t products, it is easy to see that they require $n^2(t+1) + n$ digit multiplications, which is precisely the complexity of standard lazy reduction when the products are done in schoolbook or Comba-style. This means that lazy reduction in conjunction with a subquadratic multiplication like Karatsuba-schoolbook or Karatsuba-Comba (KCM) [51] is theoretically cheaper in terms of multiplications.

Nevertheless, we argue that the new method can achieve a superior performance in practice for small and medium-size primes since it enables streamlined implementations with much less friction with memory. That is, the computation is carried out over a reduced number of registers, greatly minimizing the memory accesses. In software, the schoolbook variants are particularly efficient due to the availability of carry-preserving instructions, which precisely contribute further to reduce memory use.

To see this, let's run a comparative analysis with one of the most promising variants for software platforms, namely, a merged sum of products using radix- r interleaved Montgomery multiplication in coarsely integrated form (Algorithm 2). For the remainder of this section, we focus on the case of primes with no special shape, as assumed in the description of the algorithms from the previous section. We assume $B = 1$ in the case of an \mathbb{F}_{p^2} multiplication, and set an x64 processor as the target. We compare against

⁴We use a loose definition here: a prime should be well above 500 bits long to be considered “large”, but this varies with the computer wordsize (the smaller the wordsize the lower the threshold to consider that a prime is large).

state-of-the-art implementations of multiplication over \mathbb{F}_{p^2} , which essentially use variants of Algorithm 1 [2, 39].

First, if we perform a theoretical analysis in line with, e.g., [61], which counts the number of instructions executing a multiplication (mul), addition/subtraction (add), memory load (read) and memory store (write), the cost of one \mathbb{F}_{p^2} multiplication using Algorithm 1 is approximately given by

$$\begin{aligned}
\text{cost} &= \text{cost}_{\text{line 1}} + \text{cost}_{\text{line 2}} + \dots + \text{cost}_{\text{line 10}} \\
&= 2 \times \left(3n^2/4 \text{ muls} + (3n^2 + 4n + 2) \text{ adds} + (3n^2/2 + 15n/2 + 1) \text{ reads} + (3n^2/4 + 11n/2 \right. \\
&\quad \left. + 1) \text{ writes} \right) + 2 \times \left(n \text{ adds} + 2n \text{ reads} + n \text{ writes} \right) + \left(3n^2/4 \text{ muls} + (3n^2 + 4n + 2) \text{ adds} \right. \\
&\quad \left. + (3n^2/2 + 15n/2 + 1) \text{ reads} + (3n^2/4 + 11n/2 + 1) \text{ writes} \right) \\
&\quad + 2 \times \left(2n \text{ adds} + 4n \text{ reads} + 2n \text{ writes} \right) + \left(3n \text{ adds} + 2n \text{ reads} + 5n \text{ writes} \right) \\
&\quad + 2 \times \left((n^2 + n) \text{ muls} + 4n^2 \text{ adds} + (2n^2 + 2n) \text{ reads} + n^2 \text{ writes} \right) \\
&= (17n^2/4 + 2n) \text{ muls} + (17n^2 + 21n + 6) \text{ adds} + (17n^2/2 + 83n/3 + 3) \text{ reads} + (17n^2/4 \\
&\quad + 49n/2 + 3) \text{ writes}, \tag{4}
\end{aligned}$$

where the multiplications (lines 1, 2 and 5, Alg. 1) are assumed to be computed using Karatsuba at the top level and schoolbook underneath, and the reductions (lines 9 and 10, Alg. 1) are assumed to be computed with radix- r interleaved Montgomery multiplication in operand scanning form (schoolbook).

Now, if we do a similar cost analysis for Algorithm 2, the cost for the sum of t products is given by

$$\begin{aligned}
\text{cost} &= (tn^2 + n^2 + n) \text{ muls} + (4tn^2 + 4n^2 + 2tn) \text{ adds} + (2tn^2 + 2n^2 + 2tn + 2n) \text{ reads} \\
&\quad + (tn^2 + n^2 + 2tn) \text{ writes}.
\end{aligned}$$

And thus, the cost for a full \mathbb{F}_{p^2} multiplication, consisting of two sums of products with $t = 2$, is given by

$$\text{cost} = (6n^2 + 2n) \text{ muls} + (24n^2 + 8n) \text{ adds} + (12n^2 + 12n) \text{ reads} + (6n^2 + 8n) \text{ writes}. \tag{5}$$

If we compare costs (4) and (5), standard lazy reduction appears to beat the new method solidly for almost every operation type. However, this analysis ignores key practical considerations, as we discuss below.

Let's now perform a more practical analysis based on actual implementations of the \mathbb{F}_{p^2} multiplication for primes up to 512 bits long on an x64 platform. In order to provide a generic cost formula (but without loss of precision), the costs given below are approximations of the actual operation counts. We consider the use of carry-preserving instructions like `mulx` and `adx`, which are supported by all modern Intel and AMD processors.

In this case, the cost when using the standard lazy reduction technique (Alg. 1) is approximately given by

$$\begin{aligned}
\text{cost} &= (17n^2/4 + 2n) \text{ muls} + (17n^2/2 + 55n/2 - 9) \text{ adds} + (17n^2/4 + 40n) \text{ reads} + (2n^2 \\
&\quad + 47n/2) \text{ writes}. \tag{6}
\end{aligned}$$

And the cost of the new method using Algorithm 2 is approximately given by

$$\text{cost} = (6n^2 + 2n) \text{ muls} + (12n^2 + 6n) \text{ adds} + (6n^2 + 6n) \text{ reads} + 2n \text{ writes}. \tag{7}$$

As can be seen, in practice the memory access costs are greatly reduced thanks to the use of the general purpose registers (GPRs). Likewise, the use of carry-preserving instructions reduces the number of addition instructions significantly. While this happens across both algorithms, the improvement is much more pronounced for the new method, especially in the case of memory writes. This highlights the streamlined nature of the proposed approach, which permits to eliminate many memory accesses.

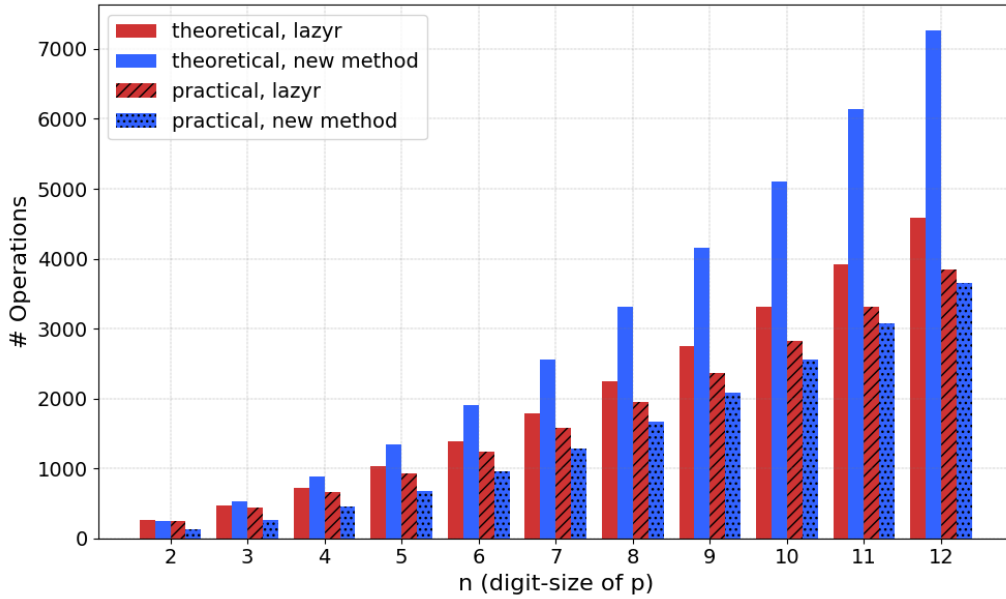


Figure 1: Comparison of instruction counts between the proposed method (Algorithm 2, $B = 1$) and the standard lazy reduction method (`lazyr`, Algorithm 1) for computing a full multiplication over \mathbb{F}_{p^2} . The counts cover all the instructions executing multiplications, additions, subtractions and memory accesses. The theoretical counts closely follow similar counts from [51, 61] and use the cost formulas (4) and (5). The practical counts, which use the cost formulas (6) and (7), are precise approximations of the actual implementation costs on an x64 platform with `mulx/adx` support for primes with different digit sizes n up to 512 bits (results for larger primes are still obtained by using the same formulas, so the measurement errors are larger in these cases, as discussed in the text).

The above analysis for primes up to 512 bits can be clearly observed in Figure 1, which displays the total number of instructions (multiplications, additions/subtractions and memory reads and writes) for different prime sizes, using the theoretical analysis and the analysis based on practical implementations. We remark that the costs for primes *greater* than 512 bits are still estimated using formulas (6) and (7). However, in practice we expect the method to become less advantageous as the prime size increases. We have verified that for larger primes the number of GPRs becomes insufficient and the implementations start to demand more storage and an increasing number of memory accesses. If we consider that from the 15 GPRs available on an x64 platform 3 are used for the input/output interface, and 3 or 4 are used to hold auxiliary values, then only 8 or 9 registers are available to hold the actual operands for the computations. Hence, it makes sense that for primes up to 512 bits, i.e., with $n \leq 8$, the proposed method achieves optimality to perform a full B -digit \times row multiplication ($B = 1$) without register spillage. This observation, conflated with the better performance of Karatsuba at larger prime sizes, makes the new method to reduce its speed advantage as the prime bitlength increases beyond 512 bits.

Experimentally, we have verified this analysis with the SIKE primes (see Appendix B). Concretely, the speed superiority of our method over standard lazy reduction starts to vanish around the 610-bit mark on an x64 platform. It is important to note, however, that the memory savings remain stable and can even slightly increase with the prime bitlength.

4 Case Study: bilinear pairings

Since the seminal papers by Sakai et al. [72] on identity-based non-interactive authentication key agreement and by Joux [56] on tripartite one-round key agreement, bilinear pairings have become a powerful tool in the design of a myriad of novel cryptographic schemes such as identity-based cryptosystems [20, 24, 29, 30] and non-interactive zero-knowledge proof systems [52, 53].

One critical drawback of pairings is their relatively expensive running time. This motivated an extraordinary effort by the research community to improve efficiency on several fronts, including the construction of pairing-friendly elliptic curves [12, 25, 49, 75], the development and improvement of the algorithms for the Miller loop and final exponentiation [11, 13, 14, 54, 77], the optimization of the explicit formulas for the curve arithmetic [3, 35, 36], and the design and optimization of tower schemes of extension fields \mathbb{F}_{p^k} [3, 18, 60]. Readers are referred to [1] for a modern take on the design and implementation of pairings.

Here, we focus on the optimization of the arithmetic over \mathbb{F}_{p^k} using the proposed algorithms. Of practical interest are extension fields of degrees 2, 4, 6, 8, 12, 24 and 48, which are commonly used to construct the tower fields in many of the popular pairings used in practice and considered for standardization [73], such as BLS12-381, BN446, BLS12-446, BN462, BLS24-509, BLS48-581 and others; see [73, Table 1] for a summary of pairing-friendly curves adopted in practice. From this set, we study the pairings with embedding degrees 12 and 24, which cover the majority of the most popular options considered for the 128- and 192-bit security levels. Accordingly, we analyze next the implementation of multiplication over \mathbb{F}_{p^2} , \mathbb{F}_{p^4} , \mathbb{F}_{p^6} , \mathbb{F}_{p^8} and $\mathbb{F}_{p^{12}}$, which are core operations in the pairing computation, and consider the following widely widespread tower scheme [1, 3, 19, 50]:

- $\mathbb{F}_{p^2} = \mathbb{F}_p[i]/(i^2 - \beta)$, with β a non-square.
- $\mathbb{F}_{p^4} = \mathbb{F}_{p^2}[s]/(s^2 - \xi)$, where $\xi = \alpha + i$ is a non-square.
- $\mathbb{F}_{p^6} = \mathbb{F}_{p^2}[v]/(v^3 - \xi)$, where $\xi = \alpha + i$ is a non-cube.
- $\mathbb{F}_{p^{12}} = \mathbb{F}_{p^4}[t]/(t^3 - s)$ or $\mathbb{F}_{p^6}[w]/(w^2 - v)$ or $\mathbb{F}_{p^2}[w]/(\omega^6 - \xi)$, where $\xi = \alpha + i$ is a non-square, non-cube.

It is important to remark that, in practice, α and β are chosen such that they have very small absolute values (typically, absolute value 1 or 2), which help improve efficiency of the extension field arithmetic.

The case of multiplication over \mathbb{F}_{p^2} . Given that for pairings, generic Montgomery multiplication (i.e., variants that do not exploit any special form in the prime) is known to provide the best performance in software, straight implementations of the variants discussed in §3 are relevant in this case. More specifically, variants that exploit the synergy between schoolbook algorithms and carry-preserving instructions are expected to outperform other approaches. Thus, we observed that the implementation of multiplication over \mathbb{F}_{p^2} can be efficiently carried out using the interleaved radix- r Montgomery multiplication variant in coarsely integrated form (Algorithm 2), with $B = 1$ to make full use of schoolbook.

The case of multiplication over \mathbb{F}_{p^4} . There are multiple choices to implement multiplication over \mathbb{F}_{p^4} . For example, it can be implemented on top of the \mathbb{F}_{p^2} arithmetic layer using our method for the multiplication over \mathbb{F}_{p^2} and Karatsuba at the \mathbb{F}_{p^4} level with or without lazy reduction. Or it could be implemented using the proposed method by seeing \mathbb{F}_{p^4} as a direct extension of \mathbb{F}_{p^2} and expressing the operations down to the base field, as discussed next.

Let $\mathbb{F}_{p^4} = \mathbb{F}_{p^2}[s]/(s^2 - \xi)$ with $\xi = \alpha + i$ as in the towering scheme above. And let $a = a_0 + a_1s$ be an element in \mathbb{F}_{p^4} , where $a_i = (a_{i,0}, a_{i,1}) \in \mathbb{F}_{p^2}$ for $i = \{0, 1\}$. Then, the multiplication $c = (c_0, c_1) = a \cdot b$ in \mathbb{F}_{p^4} can be done as follows

$$\begin{aligned}
c_{0,0} &= a_{0,0}b_{0,0} - a_{0,1}b_{0,1} + \alpha a_{1,0}b_{1,0} - \alpha a_{1,1}b_{1,1} - a_{1,0}b_{1,1} - a_{1,1}b_{1,0} \\
&= a_{0,0}b_{0,0} - a_{0,1}b_{0,1} + a_{1,0}(\alpha b_{1,0} - b_{1,1}) - a_{1,1}(b_{1,0} + \alpha b_{1,1}). \\
c_{0,1} &= a_{0,0}b_{0,1} + a_{0,1}b_{0,0} + \alpha a_{1,0}b_{1,1} + \alpha a_{1,1}b_{1,0} + a_{1,0}b_{1,0} - a_{1,1}b_{1,1} \\
&= a_{0,0}b_{0,1} + a_{0,1}b_{0,0} + a_{1,0}(b_{1,0} + \alpha b_{1,1}) + a_{1,1}(\alpha b_{1,0} - b_{1,1}). \\
c_{1,0} &= a_{0,0}b_{1,0} - a_{0,1}b_{1,1} + a_{1,0}b_{0,0} - a_{1,1}b_{0,1}. \\
c_{1,1} &= a_{0,0}b_{1,1} + a_{0,1}b_{1,0} + a_{1,0}b_{0,1} + a_{1,1}b_{0,0}.
\end{aligned} \tag{8}$$

After regrouping common coefficients and assuming that the two values $(\alpha b_{1,0} - b_{1,1})$ and $(b_{1,0} + \alpha b_{1,1})$ are pre-calculated, one can apply either Algorithm 2 or Algorithm 3 with a cost of $4 \times 4 = 16$ multiplications in the base field.

Note that, in contrast to a generic sum of products, the polynomial multiplication modulo f offers opportunities to eliminate some multiplications using Karatsuba. For example, in the term $c_{0,1}$ one could compute $a_{0,0}b_{0,1} + a_{0,1}b_{0,0}$ as $(a_{0,0} + a_{0,1})(b_{0,0} + b_{0,1}) - a_{0,0}b_{0,0} - a_{0,1}b_{0,1}$ with only *one* base field multiplication, using intermediate values from $c_{0,0}$. However, these replacements should be applied with care, since they break the algorithm's flow (recall that inner multiplications are interleaved with reduction computations) and increase memory usage, potentially neglecting any savings obtained by eliminating multiplications. Ultimately, the benefit of combining Karatsuba with the proposed algorithms might depend on the target platform (see Appendix A for details on another Karatsuba variant).

The case of multiplication over \mathbb{F}_{p^6} . Similar to the case over \mathbb{F}_{p^4} , there are multiple choices to implement multiplication over \mathbb{F}_{p^6} . For example, it can be implemented on top of the \mathbb{F}_{p^2} arithmetic layer using our method for the multiplication over \mathbb{F}_{p^2} and Karatsuba at the \mathbb{F}_{p^6} level with or without lazy reduction. Or it could be implemented using the proposed method by seeing \mathbb{F}_{p^6} as a direct extension of \mathbb{F}_{p^2} and expressing the operations down to the base field, as discussed next.

Let $\mathbb{F}_{p^6} = \mathbb{F}_{p^2}[v]/(v^3 - \xi)$ with $\xi = \alpha + i$ as in the towering scheme above. And let $a = a_0 + a_1v + a_2v^2$ be an element in \mathbb{F}_{p^6} , where $a_i = (a_{i,0}, a_{i,1}) \in \mathbb{F}_{p^2}$ for $i = \{0, 1, 2\}$. Then, the multiplication $c = (c_0, c_1, c_2) = a \cdot b$ in \mathbb{F}_{p^6} can be done as follows

$$\begin{aligned}
c_{0,0} &= a_{0,0}b_{0,0} - a_{0,1}b_{0,1} + \alpha a_{1,0}b_{2,0} - \alpha a_{1,1}b_{2,1} + \alpha a_{2,0}b_{1,0} - \alpha a_{2,1}b_{1,1} - a_{1,0}b_{2,1} - a_{1,1}b_{2,0} \\
&\quad - a_{2,0}b_{1,1} - a_{2,1}b_{1,0}. \\
&= a_{0,0}b_{0,0} - a_{0,1}b_{0,1} + a_{1,0}(\alpha b_{2,0} - b_{2,1}) - a_{1,1}(b_{2,0} + \alpha b_{2,1}) + a_{2,0}(\alpha b_{1,0} - b_{1,1}) - a_{2,1}(b_{1,0} + \alpha b_{1,1}). \\
c_{0,1} &= a_{0,0}b_{0,1} + a_{0,1}b_{0,0} + \alpha a_{1,0}b_{2,1} + \alpha a_{1,1}b_{2,0} + \alpha a_{2,0}b_{1,1} + \alpha a_{2,1}b_{1,0} + a_{1,0}b_{2,0} - a_{1,1}b_{2,1} \\
&\quad + a_{2,0}b_{1,0} - a_{2,1}b_{1,1}. \\
&= a_{0,0}b_{0,1} + a_{0,1}b_{0,0} + a_{1,0}(b_{2,0} + \alpha b_{2,1}) + a_{1,1}(\alpha b_{2,0} - b_{2,1}) + a_{2,0}(b_{1,0} + \alpha b_{1,1}) + a_{2,1}(\alpha b_{1,0} - b_{1,1}). \\
c_{1,0} &= a_{0,0}b_{1,0} - a_{0,1}b_{1,1} + a_{1,0}b_{0,0} - a_{1,1}b_{0,1} + \alpha a_{2,0}b_{2,0} - \alpha a_{2,1}b_{2,1} - a_{2,0}b_{2,1} - a_{2,1}b_{2,0}. \\
&= a_{0,0}b_{1,0} - a_{0,1}b_{1,1} + a_{1,0}b_{0,0} - a_{1,1}b_{0,1} + a_{2,0}(\alpha b_{2,0} - b_{2,1}) - a_{2,1}(\alpha b_{2,1} + b_{2,0}). \\
c_{1,1} &= a_{0,0}b_{1,1} + a_{0,1}b_{1,0} + a_{1,0}b_{0,1} + a_{1,1}b_{0,0} + \alpha a_{2,0}b_{2,1} + \alpha a_{2,1}b_{2,0} + a_{2,0}b_{2,0} - a_{2,1}b_{2,1}. \\
&= a_{0,0}b_{1,1} + a_{0,1}b_{1,0} + a_{1,0}b_{0,1} + a_{1,1}b_{0,0} + a_{2,0}(\alpha b_{2,1} + b_{2,0}) + a_{2,1}(\alpha b_{2,0} - b_{2,1}). \\
c_{2,0} &= a_{0,0}b_{2,0} - a_{0,1}b_{2,1} + a_{1,0}b_{1,0} - a_{1,1}b_{1,1} + a_{2,0}b_{0,0} - a_{2,1}b_{0,1}. \\
c_{2,1} &= a_{0,0}b_{2,1} + a_{0,1}b_{2,0} + a_{1,0}b_{1,1} + a_{1,1}b_{1,0} + a_{2,0}b_{0,1} + a_{2,1}b_{0,0}.
\end{aligned} \tag{9}$$

After regrouping common coefficients and assuming that the four values $(b_{1,0} + \alpha b_{1,1})$, $(\alpha b_{1,0} - b_{1,1})$, $(b_{2,0} + \alpha b_{2,1})$ and $(\alpha b_{2,0} - b_{2,1})$ are pre-calculated, one can apply

either Algorithm 2 or Algorithm 3 with a cost of $6 \times 6 = 36$ multiplications in the base field.

Note that similar comments to the case over \mathbb{F}_{p^4} apply to the possibility of exploiting Karatsuba.

The case of multiplication over $\mathbb{F}_{p^{12}}$. Similarly in this case, we can leverage the multiplications over \mathbb{F}_{p^2} or over \mathbb{F}_{p^6} discussed above, in combination with Karatsuba with or without lazy reduction at the $\mathbb{F}_{p^{12}}$ layer. But we can also do the computation by writing the full polynomial multiplication down to the base field level. Recall that $\mathbb{F}_{p^{12}} = \mathbb{F}_{p^2}[w]/(\omega^6 - \xi)$, where $\xi = \alpha + i$. It is straightforward to determine that, in this case, we need to compute *twelve* terms each consisting of a sum of *twelve* products (assuming the pre-calculation of ten values, similarly to multiplication over \mathbb{F}_{p^6}). Similar comments apply to the possibility of exploiting Karatsuba to products in adjacent terms.

The case of multiplication over $\mathbb{F}_{p^{24}}$ and beyond. The same methodology applies to higher extension degrees. But as we observe in the experimental results, the speedup of the method decreases with the extension degree, as reducing multiplications with Karatsuba becomes increasingly more profitable than reducing memory accesses and other small operations. However, one can expect an improvement even for extension degrees as high as 24 and 48 by combining the use of the proposed technique for the lower layers (say, for degrees 2, 4, 6 and 8) while using standard Karatsuba, with or without lazy reduction, for the top layers.

Although it has not been discussed here, the approach can be easily extended to squaring and other similar operations (e.g., sparse multiplication) over extension fields.

4.1 Performance results

In order to cover different field sizes and towerings schemes, we carry out the evaluation on two pairing-friendly curves, namely, BLS12-381 and BLS24-509, using an optimal ate pairing instantiation [77].

BLS12-381, proposed by Bowe [23], is an elliptic curve from the Barreto-Lynn-Scott (BLS) family [12] that targets the 128-bit security level and is undergoing a standardization effort in the IETF CFRG [73]. Most notoriously, this curve is widely used in zero-knowledge proofs and digital signatures in blockchain applications like Zcash [6] and Ethereum 2.0 [27]⁵. BLS12-381 is defined by the curve $E(\mathbb{F}_p) : y^2 = x^3 + 4$, with embedding degree $k = 12$. Relevant to our analysis is that, in practice, the arithmetic implementation over $\mathbb{F}_{p^{12}}$ is realized via the towerings representation: $\mathbb{F}_{p^2} = \mathbb{F}_p[i]/(i^2 + 1)$, $\mathbb{F}_{p^6} = \mathbb{F}_{p^2}[v]/(v^3 - \xi)$, and $\mathbb{F}_{p^{12}} = \mathbb{F}_{p^6}[w]/(w^2 - v)$ or $\mathbb{F}_{p^2}[w]/(\omega^6 - \xi)$, where $\xi = 1 + i$.

BLS24-509, proposed by Costello et al. [40], is also an elliptic curve from the BLS family that targets the 192-bit security level. BLS24-509 is defined by the curve $E(\mathbb{F}_p) : y^2 = x^3 + b$, with embedding degree $k = 24$. The arithmetic implementation over $\mathbb{F}_{p^{24}}$ is realized via the towerings representation: $\mathbb{F}_{p^2} = \mathbb{F}_p[i]/(i^2 + 1)$, $\mathbb{F}_{p^4} = \mathbb{F}_{p^2}[s]/(s^2 - \xi)$, $\mathbb{F}_{p^8} = \mathbb{F}_{p^4}[w]/(w^2 - s)$ or $\mathbb{F}_{p^2}[w]/(w^4 - \xi)$, and $\mathbb{F}_{p^{24}} = \mathbb{F}_{p^8}[t]/(t^3 - w)$ or $\mathbb{F}_{p^2}[t]/(t^{12} - \xi)$, where $\xi = 1 + i$.

To evaluate the proposed algorithms, we have integrated our implementations to the RELIC cryptographic library, version 0.6.0 [2]. This library contains, to our knowledge, some of the most efficient open-source implementations of pairings. In particular, it applies the generalized lazy reduction to the full extension field and elliptic curve arithmetic, as proposed in [3].

In our experiments, we use a 3.4GHz Intel Core i7-6700 (Skylake) processor with TurboBoost disabled to follow standard practice. Compilation was carried out using clang

⁵In fact, the main motivation for the design of BLS12-381 was its use for Zcash's zk-SNARK proofs.

v11.0.1 with the command `clang -O3`. Memory stack usage was obtained using `valgrind`⁶ and `massif-cherrypick`⁷.

Table 1 compares RELIC’s implementation of the extension field multiplications for BLS12-381 against the various options that we discussed for our algorithms. Likewise, Table 2 includes the analogous results for BLS24-509. We use the following notation to specify a given strategy: first, we indicate up to which layer an algorithm is applied, followed by the approach taken for the upper layers (if any). For the latter, we have two options: for the upper layers, one can use either straight Karatsuba (called “Karat”) or Karatsuba with lazy reduction (called “Karat + lazyr”). For example, if the table indicates “Alg. 2 over \mathbb{F}_{p^6} . Karat + lazyr over $\mathbb{F}_{p^{12}}$ ”, it means that we use Algorithm 2 to implement multiplication up to the \mathbb{F}_{p^6} layer, with the upper layer over $\mathbb{F}_{p^{12}}$ implemented with a formula that exploits Karatsuba with lazy reduction. In all the cases, we set $B = 1$ for Algorithm 2, which gives optimal performance on the targeted x64 platform. As noted before, this schoolbook-like algorithmic variant implementing an interleaved modular multiplication in coarsely integrated form fully exploits the availability of carry-preserving instructions. We comment that, at least on the targeted processor, the algorithm should achieve similar performance for small values of B , as long as an increase in the radix size does not put additional pressure on the register usage. For example, in our experiments, we obtained similar results for $B = 1$ and $B = 2$.

In terms of speed, Tables 1 and 2 reveal that the full use of the new method solidly beats the state-of-the-art implementations up to the \mathbb{F}_{p^6} layer in the case of BLS12-381. For the $\mathbb{F}_{p^{12}}$ multiplication, the fastest mark is achieved by using the implementation over \mathbb{F}_{p^6} and implementing the upper layer over $\mathbb{F}_{p^{12}}$ using Karatsuba. This is due to the fact that at certain threshold Karatsuba starts to outperform schoolbook algorithms when multiplications get eliminated at a sufficiently faster rate. Interestingly enough, we do not require the use of lazy reduction because a basic implementation based on Karatsuba already achieves optimal performance. This is the consequence of minimizing the cost of reduction through the proposed approach, and this greatly reduces the complexity of the implementation. Note that we also obtain a significant gain in the computation of squaring over \mathbb{F}_{p^2} . This is the result of replacing the non-interleaved Montgomery multiplication available in [2] by a faster interleaved version, given that we were not limited anymore to the old algorithmic selection that exploited lazy reduction.

In the case of BLS24-509, we experimented with a full use of the technique up to the \mathbb{F}_{p^4} level. For the \mathbb{F}_{p^8} and $\mathbb{F}_{p^{24}}$ multiplications, we use the straight implementation of the method over \mathbb{F}_{p^4} and implemented the upper layers over \mathbb{F}_{p^8} and over $\mathbb{F}_{p^{24}}$ using Karatsuba. In the latter case, the use of lazy reduction results in a slight improvement.

Another relevant point is that the experimental results confirm an important and consistent improvement in speed up to (at least) the 500-bit base field size. This can be clearly observed when comparing the speedups for the \mathbb{F}_{p^2} multiplication and squaring in the 381- and 509-bit fields. Similarly, note the increasing speedup in the multiplication over \mathbb{F}_{p^k} as k decreases for all the values of k considered in Tables 1 and 2 (this is clearly observed in the cases of full use of the method using Algorithm 2).

Additionally, we also carried out an analysis of memory usage for BLS12-381. Notably, we observe that the proposed method reduces significantly the use of memory, achieving savings in the range 43%-78% for different extension field operations and with increasing savings for higher extension degrees; see Table 1. The remarkable reduction in memory consumption is mainly due to the elimination of double-precision operations and the streamlined nature of our algorithms. Looking at the different options for \mathbb{F}_{p^6} and $\mathbb{F}_{p^{12}}$ multiplication, one can see that those that avoid lazy reduction and implement the full arithmetic using Algorithm 2 minimize the use of memory. For example, the use of

⁶<https://valgrind.org/>

⁷<https://github.com/lnishan/massif-cherrypick>

Algorithm 2 over $\mathbb{F}_{p^{12}}$ brings the most memory-friendly option for multiplication over $\mathbb{F}_{p^{12}}$, at the expense of a slight reduction in speed. We expect that similar results extend to other extension degrees.

Finally, Table 1 also reports the cycle counts for the full pairing computation using an optimal ate instantiation: on a 3.4GHz Intel Core i7-6700 Skylake machine the computation of our fastest implementation option for BLS12-381 is performed in $\sim 674 \mu\text{sec}$. As an additional data point, the same computation on a 3.2GHz Intel Core i7-8700 Coffee Lake machine is carried out in $\sim 491 \mu\text{sec}$.⁸ (compare to the 688 μsec . obtained by running the implementation from RELIC on the same platform). We remark that our implementation is much simpler and more memory-friendly, and still achieves a speedup of up to $1.40\times$ over the state-of-the-art on an x64 processor. In addition, the table also includes another implementation option in which the full $\mathbb{F}_{p^{12}}$ multiplication uses Algorithm 2. This implementation saves up to 49% in memory usage in the pairing computation while almost achieving top speed performance.

In the case of BLS24-509, Table 2 also shows a significant improvement in the running time of the full pairing. On the Skylake platform, the computation is completed in ~ 3.5 msec., which is $1.30\times$ faster than the state-of-the-art implementation from RELIC. On the Coffee Lake platform, the computation is completed in ~ 2.6 msec., which is $1.31\times$ faster than the ~ 3.4 msec. that is obtained with RELIC. As expected, when compared with the results for BLS12-381, the relatively smaller speedup obtained for BLS24-509 reflects well the use of a higher-degree tower field. Analysis for larger field sizes and higher extension degrees is left for future work.

5 Impact to Other Scenarios and Future Work

The simple but effective approach that we have proposed in this work changes the paradigm which the implementation of extension field arithmetic has long relied upon. This immediately impacts the software implementation of cryptographic schemes such as those based on bilinear pairings and supersingular isogenies. Moreover, we also expect the approach to influence the development of efficient techniques and implementations for other software platforms, constrained devices and hardware architectures, not only because of the potential speed gains but also (and maybe more critically) because of the significant savings in memory usage.

Next, we describe a few possibilities for some representative platforms, and end the section with a discussion on the impact for supersingular isogeny-based schemes.

Software platforms. In platforms with a limited number of registers there is a risk of high memory access costs. Hence, a streamlined, schoolbook-like algorithm like Algorithm 2 that minimizes memory friction and reduces the use of certain operations such as additions (e.g., when there is support for carry-preserving instructions) achieves high performance on modern x64 platforms. Alternative methods based on Karatsuba are expected to become attractive at relatively large prime sizes, when the reduction in multiplications compensates for the bumpier algorithmic flow with higher number of memory accesses and additions.

We expect a similar (if not better) situation with vectorized implementations using the recent AVX-512 vector instructions available in some Intel processors. For example, the optional extension “Integer Fused Multiply and Add” (IFMA) includes MULADD instructions that perform up to eight 52-bit multiplications followed by accumulations with 64-bit values [34]. Future work could involve studying the performance of the proposed

⁸Some fun trivia: the reported BLS12-384 implementation runs a 128-bit secure pairing in under half a millisecond, which is just slightly faster than the speed record mark hit by the BN254 pairing almost 11 years ago [3], before new attacks emerged and pushed field sizes up.

Table 1: Comparison of the speed performance (in terms of clock cycles) and memory stack usage (in terms of bytes) between the state-of-the-art implementation of an optimal ate pairing over BLS12-381 and its extension field arithmetic [2] and the optimized implementation using the method proposed in this work. The target platform is a 3.4GHz Intel Core i7-6700 (Skylake) processor.

Reference	Strategy	Speed		Memory	
		cc	%	bytes	%
\mathbb{F}_{p^2} mul					
RELIC [2]	Separated mul/rdc. Karat + lazyr	566	-	1,920	-
This work	Alg. 2 over \mathbb{F}_{p^2}	355	-37%	1,104	-43%
\mathbb{F}_{p^2} sqr					
RELIC [2]	Separated mul/rdc	451	-	1,824	-
This work	Interleaved mul/rdc	268	-41%	976	-46%
\mathbb{F}_{p^6} mul					
RELIC [2]	Separated mul/rdc. Karat + lazyr	3,376	-	6,320	-
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat over \mathbb{F}_{p^6}	2,695	-20%	2,416	-62%
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat + lazyr over \mathbb{F}_{p^6}	2,961	-12%	6,320	-0%
This work	Alg. 2 over \mathbb{F}_{p^6}	2,342	-31%	2,104	-67%
$\mathbb{F}_{p^{12}}$ mul					
RELIC [2]	Separated mul/rdc. Karat + lazyr	10,061	-	16,040	-
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat + lazyr over \mathbb{F}_{p^6} and $\mathbb{F}_{p^{12}}$	8,800	-13%	16,040	-0%
This work	Alg. 2 over \mathbb{F}_{p^6} . Karat over $\mathbb{F}_{p^{12}}$	7,858	-22%	3,928	-76%
This work	Alg. 2 over \mathbb{F}_{p^6} . Karat + lazyr over $\mathbb{F}_{p^{12}}$	8,114	-19%	13,784	-14%
This work	Alg. 2 over $\mathbb{F}_{p^{12}}$	8,315	-17%	3,544	-78%
Pairing					
RELIC [2]	Separated mul/rdc. Karat + lazyr	3.15×10^6	-	23,198	-
This work	Alg. 2 over \mathbb{F}_{p^6} . Karat over $\mathbb{F}_{p^{12}}$	2.29×10^6	-27%	12,752	-45%
This work	Alg. 2 over $\mathbb{F}_{p^{12}}$	2.30×10^6	-27%	11,792	-49%

method using the operand and product-scanning forms, in combination with different vectorization strategies.

Similar comments apply to implementations using the ARM NEON vector engine [63]. In this case, there is access to powerful, high-throughput MULADD instructions that perform up to two 32-bit multiplications followed by accumulations with 64-bit values. Thus, these instructions would favor an algorithmic variant of Eq. (3) in product-scanning form. In the case of multiplication over \mathbb{F}_{p^2} , for example, the 2-way NEON execution naturally adapts to perform the two-term computation (i.e., the operations $c_0 = a_0b_0 + a_1b_1\beta$ and $c_1 = a_0b_1 + a_1b_0$) in parallel.

For the case of scalar implementations on 64-bit ARMv8 processors, the relatively high cost of multiplication instructions might make the case for standard Karatsuba with lazy reduction. However, memory accesses are also expensive, which would favor a more streamlined execution as in the proposed algorithms. This requires actual experimentation to determine which algorithm would be optimal.

Constrained platforms. For this class of devices, the potential reduction in memory use is particularly relevant. A popular platform in this computing category is ARM Cortex-M4. For this case, one can exploit the powerful, *one-cycle* MULADD instructions available in the DSP extension. These instructions can perform a 32-bit multiplication plus 64-bit accumulation, or a 32-bit multiplication plus two 32-bit accumulations. The low cost of multiplication, added to the potential of eliminating the overhead from addition instructions in a schoolbook-like computation (e.g., see [64, §3.4]), makes our approach using Algorithm 2 a perfect fit.

Table 2: Comparison of the speed performance (in terms of clock cycles) between the state-of-the-art implementation of an optimal ate pairing over BLS24-509 and its extension field arithmetic [2] and the optimized implementation using the method proposed in this work. The target platform is a 3.4GHz Intel Core i7-6700 (Skylake) processor.

Reference	Strategy	Speed	
		cc	%
\mathbb{F}_{p^2} mul			
RELIC [2]	Separated mul/rdc. Karat + lazyr	900	-
This work	Alg. 2 over \mathbb{F}_{p^2}	554	-38%
\mathbb{F}_{p^2} sqr			
RELIC [2]	Separated mul/rdc	725	-
This work	Interleaved mul/rdc	415	-43%
\mathbb{F}_{p^4} mul			
RELIC [2]	Separated mul/rdc. Karat + lazyr	2,765	-
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat over \mathbb{F}_{p^4}	2,002	-28%
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat + lazyr over \mathbb{F}_{p^4}	2,390	-14%
This work	Alg. 2 over \mathbb{F}_{p^4}	1,885	-32%
\mathbb{F}_{p^8} mul			
RELIC [2]	Separated mul/rdc. Karat + lazyr	8,228	-
This work	Alg. 2 over \mathbb{F}_{p^2} . Karat + lazyr over \mathbb{F}_{p^4} and \mathbb{F}_{p^8}	7,081	-14%
This work	Alg. 2 over \mathbb{F}_{p^4} . Karat over \mathbb{F}_{p^8}	6,373	-23%
$\mathbb{F}_{p^{24}}$ mul			
RELIC [2]	Separated mul/rdc. Karat + lazyr	48,319	-
This work	Alg. 2 over \mathbb{F}_{p^4} . Karat over \mathbb{F}_{p^8} , Karat + lazyr over $\mathbb{F}_{p^{24}}$	41,424	-14%
Pairing			
RELIC [2]	Separated mul/rdc. Karat + lazyr	15.67×10^6	-
This work	Alg. 2 over \mathbb{F}_{p^4} . Karat + lazyr over $\mathbb{F}_{p^{24}}$	12.06×10^6	-23%

Hardware platforms. On hardware platforms like FPGAs, there are two important metrics to consider when choosing and designing an algorithm: flexibility to parallelize independent tasks, and the area-time (AT) cost. With this taken into account, our finely-integrated variant (i.e., Algorithm 3) provides an *optimal* level of parallelization at the algorithmic level without adding excessive area overhead. Note that Karatsuba and lazy reduction are not ideal in many cases because they introduce overheads and storage requirements that might hurt one or both of the hardware metrics mentioned above. Karatsuba permits to subdivide computations in multiple, smaller multiplies that can be done in parallel, but the extra circuitry can neglect a good AT trade-off on pipelined architectures. In contrast, using hardware adaptations of Algorithm 3 would naturally enable a parallel computation of up to $(t + 1)$ products (note that all the products in lines 7 and 8 in the inner for-loop are independent from each other).

Additional variants. In Appendix A, we discuss two Karatsuba variants with subquadratic complexity. These could be advantageous for applications that can afford the extra overhead and can benefit from a reduction in the number of multiplications.

We also comment that other reduction algorithms in the literature could exploit the method advantageously. For example, this is the case of Barrett reduction [16] and its interleaved version [59, §2.1], which can be extended to support merged sum of products with a unified modular reduction, as done in this work.

5.1 Other applications: supersingular isogeny-based protocols

As mentioned in §1, another notable application of the proposed technique is in the implementation of supersingular isogeny-based protocols. Recently, a series of attacks

starting with [28] rendered insecure the insignia protocols of this family, i.e., SIDH and SIKE. Since then, some works have proposed countermeasures to circumvent the attacks [47, 48, 70]. For example, in a recent work, Fouotsa, Moriya and Petit [48] propose two SIDH variants using masking. Unfortunately, the lengths of the underlying base fields of the proposed parameters are rather large and run in the thousands of bits. In this scenario, the proposed technique is expected to have a small impact (if any).

A more promising application includes the compact signature scheme SQISign [45], whose prime bitlengths belong to the ideal range of a few hundred bits. Recently, in joint work with De Feo, Leroux and Wesolowski, we adapted the proposed algorithms to a new, faster variant of SQISign [46]. Concretely, we used Algorithm 2 for p_{6983} , corresponding to a 256-bit prime of generic form, and adapted Algorithm 5 to the 254-bit prime p_{3923} , for which $p + 1 \equiv 0 \pmod{2^{64}}$. In the former case, we obtained a speedup of $\sim 1.45\times$ for the combined cost of signing and verification, while for the latter we obtained a speedup of $\sim 1.68\times$. The difference in the speedups highlights the advantage of suitably chosen primes such as p_{3923} that feature a special shape and have some room at the computer word boundaries. Future work includes the use of the algorithms for the implementation of the arithmetic over the primes for the 192- and 256-bit security levels proposed in [26]. We also comment that the recently proposed variant FastSQISignHD [41] is expected to greatly benefit from using algorithmic variants such as Algorithm 5 that are specially tailored for primes with the form $2^x \cdot 3^y - 1$, as in SIDH and SIKE.

Another interesting line of work in which our method is expected to have a significant impact corresponds to isogeny-based zero-knowledge protocols [17, 32, 44]. For example, to illustrate the impact in a recent zero-knowledge proof of isogeny knowledge [17], we modified the implementation by Basso et al. based on a 434-bit prime and observed a 23% speedup in the generation phase, and a 29% speedup in both the prove and verify stages.

References

- [1] D. F. Aranha, P. S. L. M. Barreto, P. Longa, and J. E. Ricardini. The realm of the pairings. In *Selected Areas in Cryptography – SAC 2013*, volume 8282 of *Lecture Notes in Computer Science*, pages 3–25. Springer, 2013.
- [2] D. F. Aranha, C. P. L. Gouvêa, T. Markmann, R. S. Wahby, and K. Liao. RELIC is an Efficient Library for Cryptography. <https://github.com/relic-toolkit/relic>.
- [3] D. F. Aranha, K. Karabina, P. Longa, C. H. Gebotys, and J. López. Faster explicit formulas for computing pairings over ordinary curves. In *Advances in Cryptology – EUROCRYPT 2011*, volume 6632 of *Lecture Notes in Computer Science*, pages 48–68. Springer, 2011.
- [4] R. M. Avanzi. Aspects of hyperelliptic curves over large prime fields in software implementations. In *Cryptographic Hardware and Embedded Systems – CHES 2004*, volume 3156 of *Lecture Notes in Computer Science*, pages 148–162. Springer, 2004.
- [5] R. Azarderakhsh, M. Campagna, C. Costello, L. De Feo, B. Hess, A. Hutchinson, A. Jalali, K. Karabina, D. Jao, B. Koziel, B. LaMacchia, P. Longa, M. Naehrig, G. Pereira, J. Renes, V. Soukharev, and D. Urbanik. Supersingular Isogeny Key Encapsulation (SIKE), 2017–2022. Specification available at <https://sike.org>.
- [6] E. B.-Sasson, A. Chiesa, C. Garman, M. Green, I. Miers, E. Tromer, and M. Virza. Zerocash: Decentralized anonymous payments from Bitcoin. In *IEEE Symposium on Security and Privacy – SP 2014*, pages 459–474. IEEE Computer Society, 2014.
- [7] J.-C. Bajard and S. Duquesne. Montgomery-friendly primes and applications to cryptography. *J. Cryptogr. Eng.*, 11(4):399–415, 2021.

-
- [8] S. Baktir and B. Sunar. Optimal tower fields. *IEEE Trans. Computers*, 53(10):1231–1243, 2004.
- [9] R. Barbulescu and S. Duquesne. Updating key size estimations for pairings. *J. Cryptol.*, 32(4):1298–1336, 2019.
- [10] R. Barbulescu, P. Gaudry, and T. Kleinjung. The tower number field sieve. In *Advances in Cryptology – ASIACRYPT 2015*, volume 9453 of *Lecture Notes in Computer Science*, pages 31–55. Springer, 2015.
- [11] P. S. L. M. Barreto, H. Y. Kim, B. Lynn, and M. Scott. Efficient algorithms for pairing-based cryptosystems. In *Advances in Cryptology – CRYPTO 2002*, volume 2442 of *Lecture Notes in Computer Science*, pages 354–368. Springer, 2002.
- [12] P. S. L. M. Barreto, B. Lynn, and M. Scott. Constructing elliptic curves with prescribed embedding degrees. In *Security in Communication Networks – SCN 2002*, volume 2576 of *Lecture Notes in Computer Science*, pages 257–267. Springer, 2002.
- [13] P. S. L. M. Barreto, B. Lynn, and M. Scott. On the selection of pairing-friendly groups. In *Selected Areas in Cryptography – SAC 2003*, volume 3006 of *Lecture Notes in Computer Science*, pages 17–25. Springer, 2003.
- [14] P. S. L. M. Barreto, B. Lynn, and M. Scott. Efficient implementation of pairing-based cryptosystems. *J. Cryptol.*, 17(4):321–334, 2004.
- [15] P. S. L. M. Barreto and M. Naehrig. Pairing-friendly elliptic curves of prime order. In *Selected Areas in Cryptography – SAC 2006*, volume 3897 of *Lecture Notes in Computer Science*, pages 319–331. Springer, 2006.
- [16] P. Barrett. Implementing the Rivest Shamir and Adleman public key encryption algorithm on a standard digital signal processor. In *Advances in Cryptology – CRYPTO ’86*, volume 263 of *Lecture Notes in Computer Science*, pages 311–323. Springer, 1986.
- [17] Andrea Basso, Giulio Codogni, Deirdre Connolly, Luca De Feo, Tako Boris Fouotsa, Guido Maria Lido, Travis Morrison, Lorenz Panny, Sikhar Patranabis, and Benjamin Wesolowski. Supersingular curves you can trust. In *Advances in Cryptology – EURO-CRYPTO 2023*, volume 14005 of *Lecture Notes in Computer Science*, pages 405–437. Springer, 2023.
- [18] N. Benger and M. Scott. Constructing tower extensions of finite fields for implementation of pairing-based cryptography. In *International Workshop on the Arithmetic of Finite Fields – WAIFI 2010*, volume 6087 of *Lecture Notes in Computer Science*, pages 180–195. Springer, 2010.
- [19] J.-L. Beuchat, J. E. González-Díaz, S. M., E. Okamoto, F. Rodríguez-Henríquez, and T. Teruya. High-speed software implementation of the optimal ate pairing over Barreto-Naehrig curves. In *Pairing-Based Cryptography – Pairing 2010*, volume 6487 of *Lecture Notes in Computer Science*, pages 21–39. Springer, 2010.
- [20] D. Boneh and M.K. Franklin. Identity-based encryption from the Weil pairing. In *Advances in Cryptology – CRYPTO 2001*, volume 2139 of *Lecture Notes in Computer Science*, pages 213–229. Springer, 2001.
- [21] J. W. Bos and S. Friedberger. Fast arithmetic modulo $2^X p^Y \pm 1$. In *IEEE Symposium on Computer Arithmetic – ARITH 2017*, pages 148–155. IEEE Computer Society, 2017.

- [22] J. W. Bos and S. Friedberger. Faster modular arithmetic for isogeny-based crypto on embedded devices. *J. Cryptogr. Eng.*, 10(2):97–109, 2020.
- [23] S. Bowe. BLS12-381: New zk-SNARK elliptic curve construction, 2017. <https://electriccoin.co/blog/new-snark-curve/>.
- [24] X. Boyen and B. Waters. Anonymous hierarchical identity-based encryption (without random oracles). In *Advances in Cryptology – CRYPTO 2006*, volume 4117 of *Lecture Notes in Computer Science*, pages 290–307. Springer, 2006.
- [25] F. Brezing and A. Weng. Elliptic curves suitable for pairing based cryptography. *Des. Codes Cryptogr.*, 37(1):133–141, 2005.
- [26] Giacomo Bruno, Maria Corte-Real Santos, Craig Costello, Jonathan Komada Eriksen, Michael Naehrig, Michael Meyer, and Bruno Sterner. Cryptographic smooth neighbors. Cryptology ePrint Archive, Report 2022/1439, 2022.
- [27] J. Buck. Ethereum upgrade Byzantium is live, verifies first ZK-Snark proof, 2017. <https://cointelegraph.com/news/ethereum-upgrade-byzantium-is-live-verifies-first-zk-snark-proof>.
- [28] Wouter Castryck and Thomas Decru. An efficient key recovery attack on SIDH. In *Advances in Cryptology – EUROCRYPT 2023*, volume 14008 of *Lecture Notes in Computer Science*, pages 423–447. Springer, 2023.
- [29] J. C. Cha and J. H. Cheon. An identity-based signature from gap Diffie-Hellman groups. In *Public Key Cryptography – PKC 2003*, volume 2567 of *Lecture Notes in Computer Science*, pages 18–30. Springer, 2003.
- [30] L. Chen, Z. Cheng, and N. P. Smart. Identity-based key agreement protocols from pairings. *International Journal of Information Security*, 6(4):213–241, 2007.
- [31] R. C. C. Cheung, S. Duquesne, J. Fan, N. Guillermin, I. Verbauwhede, and G. Xiaoxu Yao. FPGA implementation of pairings using residue number system and lazy reduction. In *Cryptographic Hardware and Embedded Systems – CHES 2011*, volume 6917 of *Lecture Notes in Computer Science*, pages 421–441. Springer, 2011.
- [32] Jesús-Javier Chi-Domínguez. A note on constructing SIDH-PoK-based signatures after Castryck-Decru attack. Cryptology ePrint Archive, Report 2022/1479, 2022.
- [33] P. G. Comba. Exponentiation cryptosystems on the IBM PC. *IBM Systems Journal*, 29(4):526–538, 1990.
- [34] Intel Corporation. Intel 64 and IA-32 Architectures Software Developer’s Manual. <https://www.intel.com/content/www/us/en/developer/articles/technical/intel-sdm.html>, 2021.
- [35] C. Costello, H. Hisil, C. Boyd, J. M. González Nieto, and K. Koon-Ho Wong. Faster pairings on special Weierstrass curves. In *Pairing-Based Cryptography – Pairing 2009*, volume 5671 of *Lecture Notes in Computer Science*, pages 89–101. Springer, 2009.
- [36] C. Costello, T. Lange, and M. Naehrig. Faster pairing computations on curves with high-degree twists. In *Public Key Cryptography – PKC 2010*, volume 6056 of *Lecture Notes in Computer Science*, pages 224–242. Springer, 2010.
- [37] C. Costello and P. Longa. FourQ: Four-dimensional decompositions on a \mathbb{Q} -curve over the Mersenne prime. In *Advances in Cryptology – ASIACRYPT 2015*, volume 9452 of *Lecture Notes in Computer Science*, pages 214–235. Springer, 2015.

- [38] C. Costello, P. Longa, and M. Naehrig. Efficient algorithms for supersingular isogeny Diffie-Hellman. In *Advances in Cryptology – CRYPTO 2016*, volume 9814 of *LNCS*, pages 572–601. Springer, 2016.
- [39] C. Costello, P. Longa, and M. Naehrig. SIDH Library. <https://github.com/Microsoft/PQCrypto-SIDH>, 2016–2022.
- [40] Craig Costello, Kristin E. Lauter, and Michael Naehrig. Attractive subfamilies of BLS curves for implementing high-security pairings. In *Progress in Cryptology – INDOCRYPT 2011*, volume 7107 of *Lecture Notes in Computer Science*, pages 320–342. Springer, 2011.
- [41] P. Dartois, A. Leroux, D. Robert, and B. Wesolowski. SQISignHD: New dimensions in cryptography. Cryptology ePrint Archive, Report 2023/436, 2023.
- [42] S. R. Dussé and B. S. Kaliski Jr. A cryptographic library for the Motorola DSP56000. In *Advances in Cryptology – EUROCRYPT’90*, volume 473 of *Lecture Notes in Computer Science*, pages 230–244. Springer, 1991.
- [43] Armando Faz-Hernández, Julio López Hernandez, Eduardo Ochoa-Jiménez, and Francisco Rodríguez-Henríquez. A faster software implementation of the supersingular isogeny Diffie-Hellman key exchange protocol. *IEEE Trans. Computers*, 67(11):1622–1636, 2018.
- [44] L. De Feo, Samuel Dobson, Steven D. Galbraith, and Lukas Zobernig. SIDH Proof of Knowledge. In *Advances in Cryptology – ASIACRYPT 2022*, volume 13792 of *Lecture Notes in Computer Science*, pages 310–339. Springer, 2022.
- [45] L. De Feo, D. Kohel, A. Leroux, C. Petit, and B. Wesolowski. SQISign: Compact post-quantum signatures from quaternions and isogenies. In *Advances in Cryptology – ASIACRYPT 2020*, volume 12491 of *Lecture Notes in Computer Science*, pages 64–93. Springer, 2020.
- [46] L. De Feo, A. Leroux, P. Longa, and B. Wesolowski. New algorithms for the Deuring correspondence: Towards practical and secure SQISign signatures. In *Advances in Cryptology – EUROCRYPT 2023*, volume 14008 of *Lecture Notes in Computer Science*, pages 659–690. Springer, 2023.
- [47] Tako Boris Fouotsa. SIDH with masked torsion point images. Cryptology ePrint Archive, Report 2022/1054, 2022.
- [48] Tako Boris Fouotsa, Tomoki Moriya, and Christophe Petit. M-SIDH and MD-SIDH: Countering SIDH attacks by masking information. Cryptology ePrint Archive, Report 2023/013, 2023.
- [49] D. Freeman. Constructing pairing-friendly elliptic curves with embedding degree 10. In *Algorithmic Number Theory – ANTS-VII*, volume 4076 of *Lecture Notes in Computer Science*, pages 452–465. Springer, 2006.
- [50] C. C. F. Pereira Geovandro, M. A. Simplício Jr., M. Naehrig, and P. S. L. M. Barreto. A family of implementation-friendly BN elliptic curves. *J. Syst. Softw.*, 84(8):1319–1326, 2011.
- [51] J. Großschädl, R. M. Avanzi, E. Savas, and S. Tillich. Energy-efficient software implementation of long integer modular arithmetic. In *Cryptographic Hardware and Embedded Systems – CHES 2005*, volume 3659 of *Lecture Notes in Computer Science*, pages 75–90. Springer, 2005.

- [52] J. Groth. Short pairing-based non-interactive zero-knowledge arguments. In *Advances in Cryptology – ASIACRYPT 2010*, volume 6477 of *Lecture Notes in Computer Science*, pages 321–340. Springer, 2010.
- [53] J. Groth and A. Sahai. Efficient non-interactive proof systems for bilinear groups. In *Advances in Cryptology – EUROCRYPT 2008*, volume 4965 of *Lecture Notes in Computer Science*, pages 415–432. Springer, 2008.
- [54] F. Hess, N. Smart, and F. Vercauteren. The eta pairing revisited. *IEEE Transactions on Information Theory*, 52(10):4595–4602, 2006.
- [55] D. Jao and L. De Feo. Towards quantum-resistant cryptosystems from supersingular elliptic curve isogenies. In *Post-Quantum Cryptography – PQCrypto 2011*, volume 7071 of *LNCS*. Springer, 2011.
- [56] A. Joux. A one-round protocol for tripartite Diffie-Hellman. In *Algorithm Number Theory Symposium – ANTS IV*, volume 1838 of *Lecture Notes in Computer Science*, pages 385–394. Springer, 2000.
- [57] A. Karatsuba and Y. Ofman. Multiplication of Many-Digital Numbers by Automatic Computers. *Doklady Akad. Nauk SSSR*, (145):293–294, 1962. Translation in *Physics-Doklady* 7, 595–596, 1963.
- [58] T. Kim and R. Barbulescu. Extended tower number field sieve: A new complexity for the medium prime case. In *Advances in Cryptology – CRYPTO 2016*, volume 9814 of *Lecture Notes in Computer Science*, pages 543–571. Springer, 2016.
- [59] M. Knezevic, F. Vercauteren, and I. Verbauwhede. Faster interleaved modular multiplication based on Barrett and Montgomery reduction methods. *IEEE Trans. Computers*, 59(12):1715–1721, 2010.
- [60] N. Kobitz and A. Menezes. Pairing-based cryptography at high security levels. In *International Conference on Cryptography and Coding*, volume 3796 of *Lecture Notes in Computer Science*, pages 13–36. Springer, 2005.
- [61] Ç. K. Koç, T. Acar, and B. S. Kaliski Jr. Analyzing and comparing Montgomery multiplication algorithms. *Micro, IEEE*, 16(3):26–33, 1996.
- [62] C. H. Lim and H. S. Hwang. Fast implementation of elliptic curve arithmetic in $\text{GF}(p^n)$. In *Workshop on Practice and Theory in Public Key Cryptography – PKC 2000*, volume 1751 of *Lecture Notes in Computer Science*, pages 405–421. Springer, 2000.
- [63] ARM Limited. NEON programmer’s guide, v1.0. <https://developer.arm.com/documentation/den0018/a/?lang=en>, 2013.
- [64] Z. Liu, P. Longa, G. C. C. F. Pereira, O. Reparaz, and H. Seo. FourQ on embedded devices with strong countermeasures against side-channel attacks. *IEEE Trans. Dependable Secur. Comput.*, 17(3):536–549, 2020.
- [65] Z. Liu, H. Seo, A. Castiglione, K.-K. R. Choo, and H. Kim. Memory-efficient implementation of elliptic curve cryptography for the Internet-of-Things. *IEEE Trans. Dependable Secur. Comput.*, 16(3):521–529, 2019.
- [66] P. Longa. *High-speed elliptic curve and pairing-based cryptography*. PhD thesis, University of Waterloo, 2011.

- [67] P. Longa, W. Wang, and J. Szefer. The cost to break SIKE: A comparative hardware-based analysis with AES and SHA-3. In *Advances in Cryptology – CRYPTO 2021*, volume 12827 of *Lecture Notes in Computer Science*, pages 402–431. Springer, 2021.
- [68] A. Menezes, P. Sarkar, and S. Singh. Challenges with assessing the impact of NFS advances on the security of pairing-based cryptography. In *Paradigms in Cryptology – Mycrypt 2016*, volume 10311 of *Lecture Notes in Computer Science*, pages 83–108. Springer, 2016.
- [69] P. L. Montgomery. Modular multiplication without trial division. *Mathematics of Computation*, 44(170):pp. 519–521, 1985.
- [70] Tomoki Moriya. Masked-degree SIDH. Cryptology ePrint Archive, Report 2022/1019, 2022.
- [71] N. El Mrabet and M. Joye. Guide to pairing-based cryptography. Chapman & Hall/CRC Cryptography and Network Security Series (CRC Press, 2017).
- [72] R. Sakai, K. Ohgishi, and M. Kasahara. Cryptosystems based on pairing. In *Symposium on Cryptography and Information Security – SCIS 2000, Japan*, 2000.
- [73] Y. Sakemi, T. Kobayashi, T. Saito, and R. Wahby. Pairing-Friendly Curves, 2021. <https://datatracker.ietf.org/doc/html/draft-irtf-cfrg-pairing-friendly-curves-10>.
- [74] M. Scott. Implementing cryptographic pairings. In *Pairing-Based Cryptography – Pairing 2007*, volume 4575 of *Lecture Notes in Computer Science*, pages 177–196. Springer, 2007.
- [75] M. Scott and P. S. L. M. Barreto. Generating more MNT elliptic curves. *Des. Codes Cryptogr.*, 38(2):209–217, 2006.
- [76] T. Unterluggauer and E. Wenger. Efficient pairings and ECC for embedded systems. In *Cryptographic Hardware and Embedded Systems – CHES 2014*, volume 8731 of *Lecture Notes in Computer Science*, pages 298–315. Springer, 2014.
- [77] F. Vercauteren. Optimal pairings. *IEEE Transactions on Information Theory*, 56(1):455–461, 2010.
- [78] D. Weber and T. F. Denny. The solution of McCurley’s discrete log challenge. In *Advances in Cryptology – (CRYPTO ’98)*, volume 1462 of *LNCS*, pages 458–471. Springer, 1998.

A Other algorithmic variants

In this section, we discuss two variants using Karatsuba that could be amenable for certain applications looking to reduce the number of multiplies.

The first one was already mentioned in §4 and applies to polynomial multiplication in general: for certain sums of two products, it is possible to use products from adjacent terms to convert them to a Karatsuba multiplication and save a multiplication (see the subsection “The case of multiplication over \mathbb{F}_{p^4} ”). Since the proposed algorithms interleave integer multiplications and reduction products, some care has to be taken into account to avoid excessive use of storage.

Below, we propose another option that interleaves Karatsuba multiplication with the radix- r Montgomery reduction. The new algorithm is shown in Algorithm 4 using 2-way Karatsuba.

Let's focus on the simple case with $t = 1$, i.e., a standard modular multiplication $a \cdot b \bmod p$ (the sum-of-products case easily follows). The basic idea of the algorithm is to first split operands in two halves (for a 2-way Karatsuba) using the generalized radix $r = 2^{Bw} = 2^{\lceil n/2 \rceil}$, such that operands are represented as $(a_1, a_0)_{2^{\lceil n/2 \rceil}}$. Then, from Eq. (2) and proceeding in product-scanning form, we have that:

$$\begin{aligned} u &= (a_0 b_0 + q_0 p_0) / r, \text{ with } q_0 = a_0 b_0 p' \bmod r \\ &= (u + a_1 b_0 + a_0 b_1 + q_0 p_1 + q_1 p_0) / r, \text{ with } q_1 = (u + a_1 b_0 + a_0 b_1 + q_0 p_1) p' \bmod r \\ &= u + a_1 b_1 + q_1 p_1. \end{aligned}$$

Finally, we simply replace the intermediate computation $(a_1 b_0 + a_0 b_1)$ by $(a_0 + a_1)(b_0 + b_1) - a_0 b_0 - a_1 b_1$. For B odd we proceed as other cases in this work and adjust the operations to the right digit size (lines 7 and 8 in Alg. 4).

We comment that other variants are possible and easily follow, such as for example an interleaved, 3-way Karatsuba-Montgomery multiplication that can be derived in a similar way with a splitting of operands in three parts. It is also possible to derive a SIKE-friendly version. In this case, one can conveniently set the radix to $r = 2^{e_2}$ for a prime $p = 2^{e_2} \cdot 3^{e_3} - 1$ and eliminate the multiplications by p' .

Algorithm 4 Merged sums of products using Karatsuba in a radix- r interleaved Montgomery multiplication.

Input: integers $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ s.t. $a_i, b_i \in [0, 2p)$ for $i = 0, \dots, (t-1)$ and $0 \leq \sum_{i=0}^{t-1} a_i b_i < pR$, where $R = 2^{nw}$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$, and w is the computer wordsize; the radix $r = 2^{Bw}$ s.t. $B = \lceil n/2 \rceil$, and the Montgomery constant $p' = -p^{-1} \bmod r$. Integers are represented in radix r , e.g., $a_i = (a_{i,1}, a_{i,0})_r$.

Output: the Montgomery residue $c = \sum_{i=0}^{t-1} a_i b_i \cdot R^{-1} \bmod p$ s.t. $c \in [0, 2p)$.

- 1: $u_0 \leftarrow \sum_{i=0}^{t-1} a_{i,0} \cdot b_{i,0}$
 - 2: $q_0 \leftarrow u_0 \cdot p' \bmod 2^{Bw}$
 - 3: $u \leftarrow (u_0 + q_0 \cdot p_0) / 2^{Bw}$
 - 4: $u_1 \leftarrow \sum_{i=0}^{t-1} a_{i,1} \cdot b_{i,1}$
 - 5: $u_0 \leftarrow \left(\sum_{i=0}^{t-1} (a_{i,0} + a_{i,1}) \cdot (b_{i,0} + b_{i,1}) \right) - u_0 - u_1$
 - 6: $u \leftarrow u + u_0 + q_0 \cdot p_1$
 - 7: $q_1 \leftarrow u \cdot p' \bmod 2^{(B-n \bmod 2)w}$
 - 8: $u \leftarrow (u + q_1 \cdot p_0) / 2^{(B-n \bmod 2)w}$
 - 9: $u \leftarrow u + 2^{(n \bmod 2)w} \cdot (u_1 + q_1 \cdot p_1)$
 - 10: **return** $c \leftarrow u$
-

B Effect of the prime size: analysis with the SIKE primes

In this section, we evaluate the improvements in speed performance and memory usage that can be obtained with the proposed technique for different prime sizes. For this, we use the SIKE primes [5], which cover a range of sizes of relevance for several cryptographic applications.

SIKE primes have the special form $p = 2^{e_2} \cdot 3^{e_3} - 1$, where $2^{e_2} \approx 3^{e_3}$ for integers e_2 and e_3 . For the analysis, we use Algorithm 5, which we derived from Algorithm 2 by adapting it to the special prime shape.

Cost analysis. To evaluate the performance of the proposed approach, we implemented Algorithm 5 with the Round 3 SIKE parameters p434, p503 and p610 (the number in the parameter name denotes the bitlength of the corresponding prime [5]). We also

Algorithm 5 Merged sums of products using radix- r Montgomery reduction in coarsely integrated form for a prime with the form $p = 2^{e_2} \cdot 3^{e_3} - 1$.

Input: integers $(a_0, a_1, \dots, a_{t-1})$ and $(b_0, b_1, \dots, b_{t-1})$ s.t. $a_i, b_i \in [0, 2p)$ for $i = 0, \dots, (t-1)$ and $0 \leq \sum_{i=0}^{t-1} a_i b_i < pR$, where $R = 2^{nw}$, $p = 2^{e_2} \cdot 3^{e_3} - 1$, $n = \lceil l/w \rceil$, $l = \lceil \log_2 p \rceil$, and w is the computer wordsize; $z = \lfloor e_2/w \rfloor$, $\hat{p} = (p+1)/2^{zw}$, and the radix $r = 2^{Bw}$ s.t. $B \in \mathbb{Z}$ and $0 < B \leq z$. Integers are represented in radix r , e.g., $a_i = (a_{i, \lceil n/B \rceil - 1}, \dots, a_{i,1}, a_{i,0})_r$.

Output: the Montgomery residue $c = \sum_{i=0}^{t-1} a_i b_i \cdot R^{-1} \bmod p$ s.t. $c \in [0, 2p)$.

```

1:  $u \leftarrow 0$ 
2: for  $j = 0$  to  $\lceil n/B \rceil - 1$  do
3:    $u \leftarrow u + \sum_{i=0}^{t-1} a_{i,j} \cdot b_i$ 
4:    $q \leftarrow u \bmod 2^{Bw}$ 
5:    $u \leftarrow \lfloor u/2^{Bw} \rfloor + 2^{(z-B)w} q \cdot \hat{p}$ 
6: if  $n \bmod B \neq 0$  then
7:    $u \leftarrow u + \sum_{i=0}^{t-1} a_{i, \lceil n/B \rceil - 1} \cdot b_i$ 
8:    $q \leftarrow u \bmod 2^{(n \bmod B)w}$ 
9:    $u \leftarrow \lfloor u/2^{(n \bmod B)w} \rfloor + 2^{(z-n \bmod B)w} q \cdot \hat{p}$ 
10: return  $c \leftarrow u$ 

```

evaluate the SIKE prime p377 proposed by Longa et al. [67]. We compare against the performance of the state-of-the-art implementations of the same SIKE primes from the SIDH library v3.4 [39] and from [67]. These libraries implement the integer multiplication in two layers using Karatsuba (upper layer) and schoolbook (lower layer). The reduction part is implemented using a non-interleaved radix- r Montgomery reduction with $B > 1$, specialized to SIDH/SIKE primes (see [43, Alg. 6]). As is standard, these libraries use lazy reduction for the multiplication over \mathbb{F}_{p^2} , following Algorithm 1.

Table 3 presents the performance comparison (in terms of clock cycles) for the \mathbb{F}_{p^2} multiplication on an x64 processor, specifically, a 3.4GHz Intel Core i7-6700 (Skylake) processor. All the implementations in the comparison are written in assembly language, and were compiled and tested on the same platform using clang v6.0.1 with the command `clang -O3`. The table also includes a detailed instruction count of all the implementations, including multiplications, additions, subtractions and other instructions. The columns “read” and “write” present counts of all the corresponding instructions that require a memory access operation. We remark that the total instruction counts are provided as an additional data point only and should not be considered to follow actual performance with high-precision. Especially in the case of the targeted platform, its superscalar, deeply pipelined microarchitecture makes extremely difficult to extract performance data from a straight instruction count. Nevertheless, it can provide relevant information for a first-order comparison of the different algorithms.

Firstly, we observe that the proposed method achieves much better speed performance in all the cases, even though it requires a higher number of multiplication instructions. This is due to the significant reduction of other operations, especially of those requiring read/write memory accesses. Another relevant aspect is that at certain threshold the operand sizes become too large and the lack of enough general purpose registers forces the use of many more memory access instructions. This can be observed for the largest prime under analysis, i.e., p610, which precisely returns the lowest speedup. In the rest of the cases, the speedup goes from $\sim 1.17\times$ up to $1.31\times$, with the speedup increasing as the size of the prime decreases. This is consistent with the results from existing literature that show that Karatsuba becomes more profitable as sizes grow (see §3). Nevertheless, we demonstrate that, for the case of a quadratic extension field, schoolbook can still be

Table 3: Comparison of instruction counts, speed performance (in terms of clock cycles) and memory stack usage (in terms of bytes) between the proposed method (Algorithm 5, $B = 1$) and the state-of-the-art implementations of multiplication over \mathbb{F}_{p^2} for the SIDH and SIKE protocols [5, 39, 67]. The target platform is a 3.4GHz Intel Core i7-6700 (Skylake) processor. The comparison covers the Round 3 SIKE primes p434, p503 and p610, and also the prime p377 proposed in [67]. The instruction counts cover all the instructions executing multiplications, additions, subtractions and memory accesses. The column “others” includes any other additional instructions such as logical and shift instructions.

Reference	Instruction count							Speed		Memory	
	read	write	mul	add	others	total	%	cc	%	bytes	%
p377											
[67]	360	146	117	438	190	1,251	-	352	-	1,096	-
This work	227	18	192	412	68	917	-26.7%	269	-23.6%	712	-35.0%
p434											
[39]	492	196	179	589	164	1,620	-	440	-	1,224	-
This work	292	21	252	537	74	1,176	-27.4%	341	-22.5%	784	-35.9%
p503											
[39]	568	225	208	677	236	1,914	-	514	-	1,320	-
This work	366	24	336	709	90	1,525	-20.3%	440	-14.4%	832	-37.0%
p610											
[39]	903	368	345	1,019	182	2,817	-	762	-	1,536	-
This work	715	170	500	1,090	117	2,592	-8.0%	734	-3.7%	952	-38.0%

much faster for primes up to around 500 bits⁹.

Table 3 also summarizes the stack memory usage corresponding to each parameter set. The figures were obtained using `valgrind` and `massif-cherrypick`, as in §4. As can be seen, our approach achieves a significant reduction in memory consumption that is well above 35%. This is obtained consistently across the different parameter sets, with even a slight increase in the savings as the prime size goes up. This is consistent with the memory analysis in §4, although for pairings the savings are even greater for higher degree extension fields when the use of the tower-based approach is minimized.

⁹In the case of generic primes, as used for pairings, we observed a consistent speedup up to the 500-bit mark; see §4.1.