

Accelerating Homomorphic Computations on Rational Numbers

Angela Jäschke^(✉) and Frederik Armknecht

University of Mannheim, Mannheim, Germany
{jaeschke,armknecht}@uni-mannheim.de

Abstract. Fully Homomorphic Encryption (FHE) schemes are conceptually very powerful tools for outsourcing computations on confidential data. However, experience shows that FHE-based solutions are not sufficiently efficient for practical applications yet. Hence, there is a huge interest in improving the performance of applying FHE to concrete use cases. What has been mainly overlooked so far is that not only the FHE schemes themselves contribute to the slowdown, but also the choice of data encoding. While FHE schemes usually allow for homomorphic executions of algebraic operations over finite fields (often \mathbb{Z}_2), many applications call for different algebraic structures like signed rational numbers. Thus, before an FHE scheme can be used at all, the data needs to be mapped into the structure supported by the FHE scheme.

We show that the choice of the encoding can already incur a significant slowdown of the overall process, which is independent of the efficiency of the employed FHE scheme. We compare different methods for representing signed rational numbers and investigate their impact on the effort needed for processing encrypted values. In addition to forming a new encoding technique which is superior under some circumstances, we also present further techniques to speed up computations on encrypted data under certain conditions, each of independent interest. We confirm our results by experiments.

Keywords: Confidential machine learning · Fully homomorphic encryption · Encoding · Implementation

1 Introduction

Fully Homomorphic Encryption (FHE) is a very promising field of research because it allows arbitrary computations on encrypted data. This means that data can be outsourced securely without sacrificing functionality, as any operation one would like to perform on the data can also be performed on the encrypted data by a third party without divulging information. With a powerful enough encryption scheme, this third party may even apply its own proprietary algorithm, like a machine learning algorithm, to the encrypted data such that the result divulges nothing about the algorithm that was applied - this is the setting we will assume. While multiparty computation also offers this kind of

confidential computation, it requires frequent interaction between the involved parties, which seems unfortunate for the goal of outsourcing computation. For this reason, we instead focus on FHE, which allows a non-interactive solution. Unfortunately, FHE-based solutions today are still very slow and thus not very practical. Since a ciphertext can become undecryptable if too many consecutive multiplications are computed, multiplicative depth is often key in FHE computations. In so-called leveled FHE schemes, one can adjust the encryption scheme to support a predetermined multiplicative depth, where the scheme becomes slower the larger the depth is. Thus, minimizing depth is one of our goals in this paper. Another approach for handling the problems that come with consecutive multiplications, which we opted for because of very large depths in our use cases, is called bootstrapping. Here, the ciphertext is “cleaned up” after multiplication, but this operation takes very long and constitutes the bottleneck when used. Hence, minimizing the total number of multiplications is another of our goals.

Because of these efficiency problems, there is currently much research on improving the efficiency of the schemes themselves on the one hand, and on designing algorithms that are particularly suited to FHE, i.e., through minimal multiplicative depth, on the other hand. While this is certainly a valuable contribution for some use cases, we feel that in general the algorithms one wants to perform on the data are predetermined and not up for discussion. At first glance, this might seem to imply that there is little potential for improvement apart from improving the schemes themselves, but we show that this is indeed not the case.

Generally, suppose one has an FHE scheme $\mathcal{E} = (\text{Gen}, \text{Enc}, \text{Dec})$ with plaintext space \mathcal{M} and ciphertext space \mathcal{C} , and there is a function $g : \mathcal{M}^z \rightarrow \mathcal{M}$ for some $z \in \mathbb{N}$. Then a Fully Homomorphic Encryption scheme promises that there exists a corresponding function $g^* : \mathcal{C}^z \rightarrow \mathcal{C}$ with

$$\text{Dec}(\text{sk}, g^*(\text{Enc}(\text{pk}, m_1), \dots, \text{Enc}(\text{pk}, m_z))) = g(m_1, \dots, m_z).$$

However, plaintext spaces for encryption schemes are usually some finite field $GF(p^d)$ for some prime p and power d , so if we want to work with elements from a different structure S (like the rational numbers), we must first map them¹ to the plaintext space using an encoding $\pi : S \rightarrow \mathcal{M}^k$ and then perform a function on the plaintext values that emulates the function on S . For a better understanding, suppose we have an encryption scheme like above. Then, if we want to evaluate a function $f : S^n \rightarrow S$ on encrypted data, we must first turn f into a function $g : (\mathcal{M}^k)^n \rightarrow \mathcal{M}^k$ on the plaintext space (where \mathcal{M}^k emulates S) and then execute the function $g^* : (\mathcal{C}^k)^n \rightarrow \mathcal{C}^k$ that corresponds to g . This is illustrated in Fig. 1.

As it turns out, there is often no unique function g for a given function f , but instead several different ones which depend on the chosen encoding function π . This also means that the most we can aim for in terms of efficiency in evaluating a function f on encrypted data is not f itself, but rather its emulation g on

¹ For example, if $S = \{x \in \mathbb{Z} \mid 0 \leq x \leq 7\}$ (i.e., numbers representable by 3 bits) but the plaintext space of the encryption scheme is only $\mathcal{M} = \{0, 1\}$, we could map $\pi : S \rightarrow \mathcal{M}^3$.

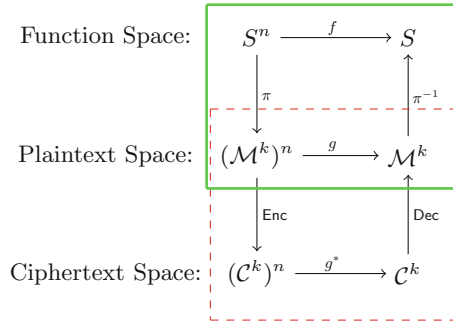


Fig. 1. Steps in homomorphic evaluation

the plaintext space. As it turns out, the increase here is not negligible: While the Perceptron, which we evaluate in Sect. 6.3 on encrypted data, runs almost instantaneously (roughly 0.004 s) for ten rounds when computing on unencrypted rational numbers, the evaluation of the same algorithm emulated on the plaintext space (i.e., still unencrypted) takes over 120 s for the same parameters even with our most efficient encoding in the plaintext space. This shows that though largely ignored until now, the overhead that comes from switching from the function f to g can be substantial and must equally be addressed to make FHE applications as efficient as possible. Thus, while previous work on making computations with FHE more efficient has focused primarily on the area inside the dashed red rectangle in Fig. 1, we investigate how to improve efficiency through the right choice of π and subsequently g , represented by the solid green rectangle. Motivated by the idea of outsourcing actual data and running existing algorithms on it, we face the challenges of encoding rational numbers (as opposed to elements of finite fields or unsigned integers) and of incorporating basic operations like addition, multiplication and comparison, which are needed for many popular algorithms.

1.1 Our Contribution

We address the above challenges and try to minimize total number of multiplications (and the multiplicative depth) of g through appropriate choices in π . We also examine some further optimizations which increase efficiency under certain assumptions and are of independent interest. As a concrete application, we apply our results to two use cases from machine learning, the Perceptron and the Linear Means Classifier, and see that the right choice of π can make a significant difference in terms of multiplicative depth, total number of multiplications, and in terms of runtime, for which we encrypted the data with the HELib library. To this end:

- We present a new method for working with encrypted rational numbers by solving the problem that the number of digits of precision doubles with each multiplication. We show how to remove the extra digits and bring the number

back down to a predefined precision level, greatly improving performance without leaking information about the function that was applied.

- We investigate two different popular encodings with regard to efficiency in emulating basic operations on rational numbers like comparison, addition and multiplication, and present a hybrid encoding that surpasses the two traditional ones both in theory (as measured by total bit additions, multiplications and required multiplicative depth) and in terms of actual runtime for large sizes.
- We the comparison of two encrypted numbers and present an easier way for comparing numbers to 0 which takes almost no time.
- We show how to increase efficiency in the case that the numbers are bounded, like in real-world applications where values lie in some known range.
- We confirm our results by implementing the Perceptron, a fundamental algorithm in machine learning, and running it using the different encodings, as well as a polynomial like that used for Linear Means Classification.

As a quick preview, consider Fig. 2, which shows theoretical bounds on the number of bitwise additions and multiplications as well as extrapolated runtime needed to apply a Linear Means Classifier with each of the three encodings for different numbers of features. We can see that our new hybrid encoding mechanism is superior in all three aspects, making it an attractive choice.

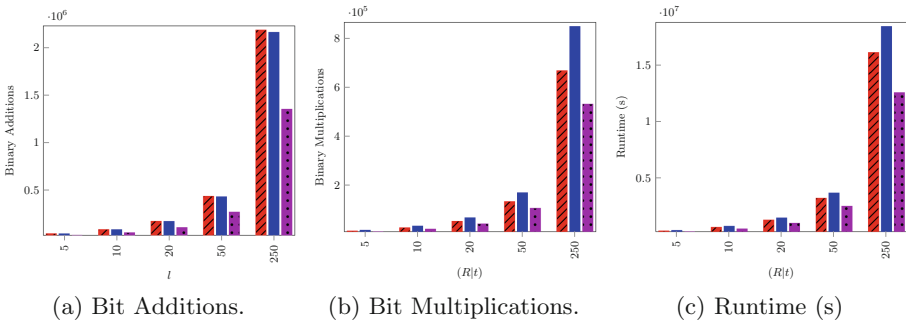


Fig. 2. Bounds for the number of bitwise additions and multiplications as well as runtime for evaluating Linear Means Classifier with l features of length 30 for different l using Two's Complement \bullet (lines), Sign-Magnitude \bullet (solid) and Hybrid Encoding \bullet (dotted) (Color figure online)

1.2 Outline

We start by giving an overview of related work in Sect. 2. In Sect. 3, we give some background on Fully Homomorphic Encryption and the challenges faced when working with rational numbers, as well as on the two encodings we use. In Sect. 4, we show how to emulate the addition, multiplication and comparison of encoded numbers using just binary additions and multiplications and analyze

complexity. Section 5 presents different ways of accelerating computations on encrypted data, and Sect. 6 gives some motivation and necessary background on machine learning before using two algorithms from this field to demonstrate the effects of our improvements. Lastly, Sect. 7 gives our conclusion and an insight into future work.

2 Related Work

While encryption schemes that allow one type of operation on ciphertexts are well understood and have a comprehensive security characterization [4], Fully Homomorphic Encryption, which allows both unlimited additions and multiplications, was only first solved in [19]. Since then, numerous other schemes have been developed, for example [9, 10, 13, 14, 16, 21, 26]. An overview can be found in [3]. There have been several works concerning actual implementation of FHE, like [20] (homomorphically evaluating the AES circuit), [7] (predictive analysis on encrypted medical data), or [22] (machine learning on encrypted data), and there are two publicly available libraries [1, 18]. [24] discusses whether FHE will ever be practical and gives a number of possible applications, including encrypted machine learning. Most recently, two publications regarding encoding rational numbers for FHE have appeared, illustrating what an important topic this is: [12] examines encoding rational numbers through continued fractions (restricted to positive rationals and evaluating linear multivariate polynomials), whereas [15] focuses on most efficiently embedding the computation into a single large plaintext space. Another work that explores similar ideas as [15] and also offers an implementation is [17].

While the idea of being able to privately evaluate machine learning algorithms is certainly intriguing, the overwhelming majority of work in this area considers multiparty computation, which requires interaction between the client and the server during computation and is thus a different model. Examples include [8, 25, 28], and works like [23, 27] concern themselves with efficiency measures and circuit optimizations specific to multiparty computation. Another line of research regarding confidential machine learning, e.g. [7] and again [8], focuses on a scenario where the model being computed and/or evaluated is publicly known - a scenario we explicitly exclude. Other work like [11] restricts itself to unsigned integers, making all involved circuits much less complex. Work like [5] considers recommender systems, but in a scenario which becomes insecure if too many fresh encryptions are available. Closest to our work is [22], which restricts itself to machine learning algorithms like the Linear Means Classifier and Fishers Linear Discriminant Classifier, which can be expressed as polynomials of low degree, and focuses on the classification, not the derivation of the model. Their encoding of input data is also restricted to functions with few multiplications.

We stress that until now, all approaches dealing with rational numbers either restrict computations to positive integers, or the multiplicative depth of the computation must be known beforehand. Our approach is the first to actually tackle the problem of computing on rational numbers with no further assumptions, and offers other improvements if some assumptions can be made.

3 Background

3.1 FHE and Efficiency Metrics

Fully Homomorphic Encryption (FHE) describes a class of encryption schemes that allow arbitrary operations on encrypted data. This would, in theory, enable outsourcing of encrypted data to an untrusted cloud service provider, who could still perform any operations the user wishes. This means that we can protect privacy (as opposed to uploading the data in unencrypted form) while maintaining functionality (as opposed to uploading data encrypted under conventional schemes). Unfortunately, FHE today it is still rather slow, although huge advancements have been made in the last six years.

Because of this, one of our measures for efficiency is the number of bit additions and multiplications performed, as this would translate directly into the number of homomorphic additions and multiplications performed if the data were encrypted. Note that in schemes today, homomorphic multiplication tends to be computationally more expensive than addition.

In our analysis of computational effort, we also include the multiplicative depth: Many publications today use *Leveled Fully Homomorphic Encryption*, which is related to Fully Homomorphic Encryption in that arbitrary functions f can be performed on the encrypted data, but the multiplicative depth of f must be known beforehand, and efficiency of the encryption scheme decreases as this number increases. Multiplicative depth measures how many consecutive multiplications are performed. For example, the polynomial $x_1 \cdot x_2 + x_1 \cdot x_3 + x_2 \cdot x_3$ has 3 multiplications in total, but a multiplicative depth of only 1. These leveled schemes can be more efficient than pure FHE schemes for small depths, but if more than the allowed number of consecutive multiplications are performed, decryption may return the wrong result. To this end, we include multiplicative depth in our analysis and aim to minimize it as one of our goals. We would, however, like to point out that if one uses bootstrapping, as we did in our implementations, depth becomes less of an issue and the total number of multiplications is the main factor determining runtime.

3.2 From Unsigned Integers to Rationals of Arbitrary Precision

In previous work (e.g. [6], see also Sect. 2), rational numbers have often been approximated by multiplying with a power of 10 and rounding, but note that when multiplying two rational numbers with k bits of precision, we obtain a number with $2k$ bits of precision (whereas addition does not change the precision). If we are working on unencrypted numbers, we might just round to obtain k bits of precision again, or we could truncate (truncation after k bits yields the same accuracy as rounding to $k - 1$ bits). However, things become more difficult if we will be operating on encrypted data, as rounding is generally not possible here and thus these extra bits of precision accumulate. To see this, suppose a precision of k digits is required. One would usually multiply the rational number with 10^k and round (or truncate) to the nearest integer, which is then encoded

and encrypted. Dividing the decrypted decoded number by 10^k again yields the rounded rational. However, the problem of doubling precision with multiplication is prevalent here. Consider what would happen if we were to multiply two such numbers: Suppose we have two rational numbers a and b that we would like to encode as integers a' and b' with k digits of precision, so we get $a' = a \cdot 10^k$ and $b' = b \cdot 10^k$ (rounded to the nearest integer). Multiplying a' and b' , we get $c'' = a' \cdot b' = a \cdot 10^k \cdot b \cdot 10^k = (a \cdot b) \cdot 10^{2k}$. Thus, having reversed the encoding, the obtained value c'' must be divided by 10^{2k} . This is a problem because we cannot remove the extra bits by dividing by 10^k , so the party performing the algorithm must now divulge what power of 10 to divide the obtained result by. This leaks information about the multiplicative depth of the function used and thus constitutes a privacy breach for the computing party. Additionally, there is also the problem during computation that the sizes of the encoded numbers will increase substantially.

To solve this problem, we propose the following approach: Instead of scaling by a power of 10, we multiply by a power of 2 and truncate to obtain an integer that we will encode in binary fashion, so that we can later encrypt each bit separately. This eliminates the above problem: Multiplying two numbers a' and b' with k bits of precision still yields $c'' = (a \cdot b) \cdot 2^{2k}$, but since we are encoding bit by bit, dividing by 2^k and truncating corresponds to merely deleting the last k (encrypted) bits of the product. Thus, the party performing the computations can bring the product c'' back down to the required precision after every step by discarding the last k bits and thus obtaining $c' = a \cdot b \cdot 2^k$, meaning that the party which holds the data must always divide the decoded result by 2^k no matter what operations were applied. This has the benefit of not only hiding the data from the computing party, but also hiding the function from the party with the data.

3.3 Two's Complement

Having determined that we will be encoding bit for bit to support arbitrary precision without information leakage, we must now decide on how exactly we want to represent a rational number (which has been scaled to be a signed integer). For unsigned integers, binary representation is well known: Given an integer $a \geq 0$, we write it as $a = \sum_{i=0}^n a_i \cdot 2^i$ where $n = \lfloor \log_2(|a|) \rfloor$ and $a_i \in \{0, 1\}$ to obtain a $n + 1$ -bit string $a_n a_{n-1} \dots a_1 a_0$.

To incorporate negative numbers, the most popular encoding is called *Two's Complement*: Here, we write an integer a as $a = a_{n+1} \cdot (-2^{n+1}) + \sum_{i=0}^n a_i \cdot 2^i$ where $n = \lfloor \log_2(|a|) \rfloor$ and $a_i \in \{0, 1\}$. This means that the most significant bit (MSB) encodes the negative value -2^{n+1} and is thus 1 exactly when $a < 0$. As an example, consider the bitstring 1011, which encodes $1 \cdot (-2^3) + 0 \cdot 2^2 + 1 \cdot 2 + 1 \cdot 1 = -8 + 2 + 1 = -5$.

The most notable aspect for Two's Complement is that for multiplication to work, the inputs must first be encoded as numbers of the length that the output

will have, i.e., when multiplying numbers of lengths n and m , both inputs lengths need to be increased to $n + m$ before multiplication. This procedure, called *sign extension*, is done by replacing the first bit with the appropriate number of copies of itself. In the above example, if we needed to extend the 4-bit string 1011 to length 8, it would result in 11111011, which still encodes -5 .

3.4 Sign-Magnitude

While Two's Complement may be the most popular encoding of signed integers, it is not the only one: *Sign-Magnitude* encoding formalizes the most intuitive idea of having an extra bit that determines the sign. Conventionally, this is the most significant bit, which is 1 when a number is negative and 0 when a number is positive. Thus, for example, the number $5 = 0101$ and $-5 = 1101$. This notation suffers from the fact that there are two encodings of 0 ($0 = 00\dots00$ and $-0 = 10\dots00$) and is seldom used, but we will later see how this slightly unconventional encoding can help us.

We would like to point out that addition in this encoding is much more involved than in Two's Complement: Here, we need to add the absolute values and keep the sign bit if both inputs have equal signs, and otherwise compare the two inputs, subtract the smaller from the larger absolute value, and keep the sign of the input with the larger absolute value. Obviously, expressing this operation as a polynomial is considerably more involved than the straightforward addition used in Two's Complement. However, in multiplication, Sign-Magnitude encoding does not need sign extension, and addition of the rows in multiplication can use the straightforward addition instead of the above one, so this problem does not carry over to multiplication.

4 Basic Operations and Their Performance

Having introduced two different ways of encoding, this section will now examine both the theoretical complexity and actual performance of elementary operations. All computations were done on a virtual machine with 5 GB of RAM running Ubuntu 14.04 LTS (running on a Lenovo Yoga 2 Pro with a Intel i7-4500U processor with 1.8 GHz and 8 GB of RAM with Windows 8.1). We give the number of binary additions and multiplications as well as multiplicative depth required for these elementary operations. Due to space limitations, we omit how these values were determined, but we used straightforward methods to turn the functions into polynomials over $\{0, 1\}$ and derived the number of bit additions and multiplications as well as the multiplicative depth. We note that we also implemented all our functions with a subroutine that counts these values to ensure that the formulas are correct. Runtimes were obtained for values encrypted with the HELib library [1].

4.1 Note on Comparisons

As already mentioned, Sign-Magnitude uses a comparison in its addition function, making the comparison function an important building block. We note, however, that when comparing a number with 0, there is an easier way (see Sect. 5.2). For the general case (and used in Sign-Magnitude’s addition procedure), the effort of comparing two arbitrary numbers is:

- | Two’s Complement: | Sign-Magnitude: |
|----------------------------------|--------------------------------------|
| • $3n$ binary additions | • $10n - 3$ binary additions |
| • $n + 1$ binary multiplications | • $6n - 2$ binary multiplications |
| • a multiplicative depth of n | • a multiplicative depth of $2n - 1$ |

We can see that Two’s Complement is more efficient for comparing encrypted numbers.

4.2 Addition

We will now compare addition of two n -bit numbers for Two’s Complement and Sign-Magnitude encoding. The computational effort is:

- | Two’s Complement: | Sign-Magnitude: |
|---------------------------------|--------------------------------------|
| • $5n - 2$ binary additions | • $73n - 17$ binary additions |
| • n binary multiplications | • $28n + 4$ binary multiplications |
| • a multiplicative depth of n | • a multiplicative depth of $2n + 2$ |

As we can see, Two’s Complement again does better in theory. In practice (i.e., counted by our program), we get as values the number of operations and runtime as shown in Fig. 3. These diagrams show that Two’s Complement is indeed superior to Sign-Magnitude where addition is concerned.

4.3 Multiplication

In this section, we will examine the multiplication of an n -bit number with an m -bit number. Heuristically, we expect Sign-Magnitude to do better here: Instead of the costly “normal” Sign-Magnitude addition operation which uses a comparison circuit, we can use regular textbook binary addition to add up the rows encountered in multiplication, so the fact that addition of two n -bit Sign-Magnitude numbers is much more expensive than that of two n -bit Two’s Complement numbers does not weigh in here. On the other hand, because of the sign extension necessary in Two’s Complement multiplication, not only are the rows longer ($n + m$ as compared to n), but there are also more of them ($n + m$ as opposed to m), so we must do more additions of longer bitstrings. We examine the effort required:

Two’s Complement:

- $\frac{5(m^2+n^2)-19(m+n)}{2} + 5mn + 10$ binary additions
- $\frac{(m+n-3)(m+n)}{2} + mn + 1$ binary multiplications
- a multiplicative depth of $\lceil \log_2(m+n) \rceil \cdot (m+n-1) - 2^{\lceil \log_2(m+n) \rceil} + 2$

Sign-Magnitude: Due to changing intermediate lengths during row additions (which depend on both n and m instead of just $n+m$ as in Two’s Complement), an exact formula would be very involved and hardly informative. Thus, we present a formula for an upper bound which already shows that SM is superior to TC for multiplication. To this end, we now have two data sets for Sign-Magnitude in the diagrams 3b, d and f in Fig. 3 regarding the number of operations: One shows the exact numbers as counted by an instruction in our program (and verified manually), and one shows the bounds as given by the following formulas:

- $(2^{\lceil \log_2(m-1) \rceil} - 1) \cdot (5n - 7) + (2^{\lceil \log_2(m-1) \rceil - 1} - 1) \cdot 5 \cdot \lceil \log_2(m-1) \rceil$
binary additions at most
- $(n-1) \cdot (m-1) + (2^{\lceil \log_2(m-1) \rceil} - 1) \cdot (n-1) + (2^{\lceil \log_2(m-1) \rceil - 1} - 1) \cdot \lceil \log_2(m-1) \rceil$
binary multiplications at most
- A multiplicative depth of at most $\frac{1}{2} \lceil \log_2(m-1) \rceil \cdot (\lceil \log_2(m-1) \rceil + 2n - 5) + 2^{\lceil \log_2(m-1) \rceil}$

Concrete values and runtimes can be seen in Fig. 3 and as we can see, Two’s Complement performs much worse, as expected. *Thus, Two’s Complement encoding is superior for addition and comparison, but inferior for multiplication.*

5 Accelerating Computations

In this section, we will discuss several optimizations to make computations on encrypted data more efficient.

5.1 Hybrid Encoding

Since we have seen in the previous sections that Two’s Complement encoding always performs better than Sign-Magnitude except for multiplication (where it is much worse), we propose the following approach, called Hybrid Encoding: We work with Two’s Complement encoding, but when we want to multiply, we convert the numbers to their representations in Sign-Magnitude, perform the multiplication there, and convert the result back. As we will see, this is indeed more efficient than regular Two’s Complement multiplication. To do this, we must first determine how to convert numbers from their representation in Two’s Complement to their Sign-Magnitude form and vice versa, so suppose we have a number a under one encoding α (either Two’s Complement or Sign-Magnitude), denoted a_α , and wish to transform it into its representation under the other encoding β , denoted a_β . For numbers with MSB 0, both encodings are actually the same ($a_\alpha = a_\beta$), so in this case we do nothing. If the number has a MSB of

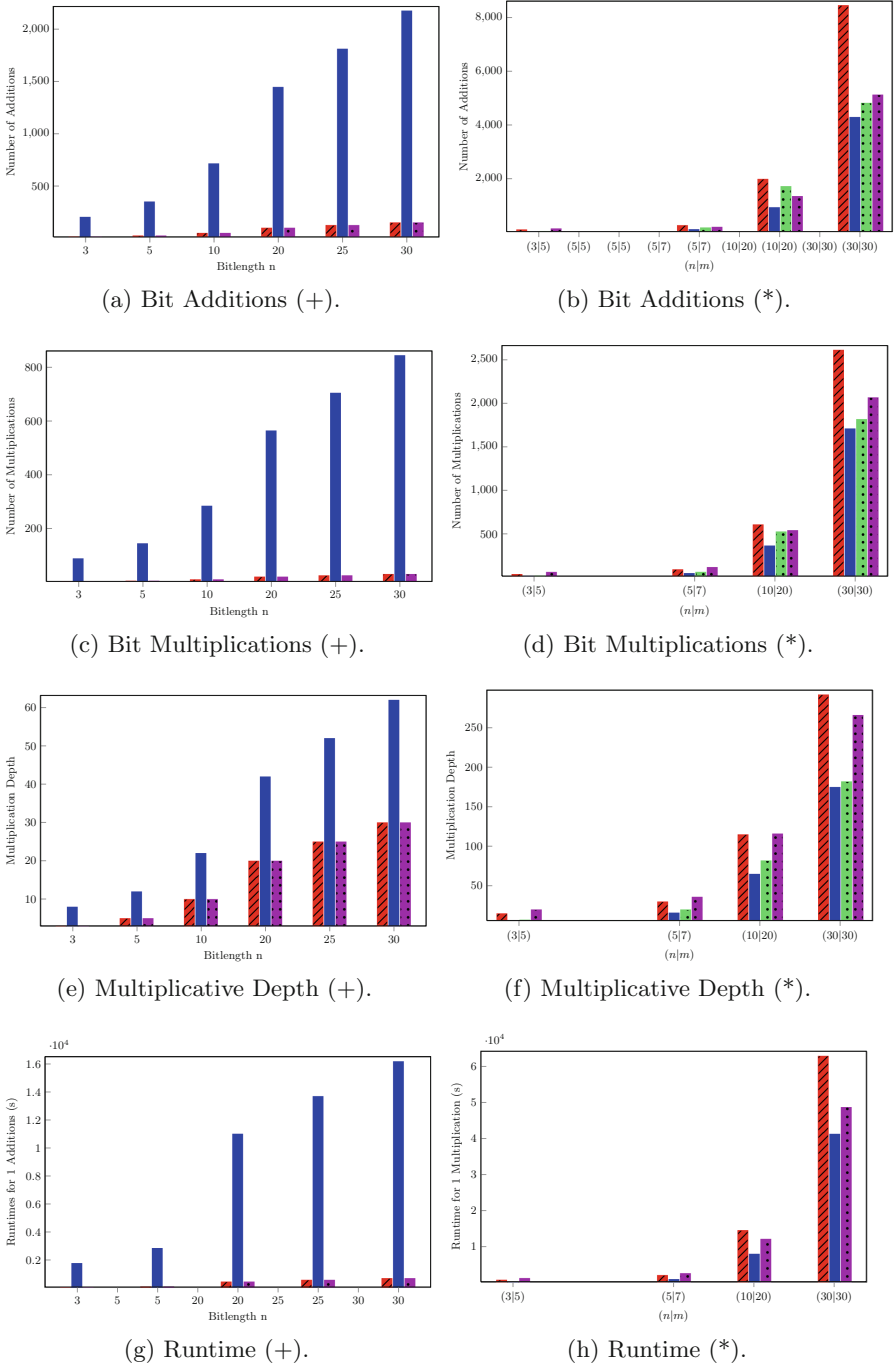


Fig. 3. Comparison of addition (+) and multiplication (*) for Two's Complement (red lines), exact values for Sign-Magnitude (counted by program) (blue solid), upper bound for Sign-Magnitude for multiplication (green dotted) and our new Hybrid Encoding (purple dotted). Runtimes for data encrypted with HELib (Color figure online).

1, we compute its negation ($a_\alpha \mapsto -a_\alpha$), which is the same for both encodings as it has MSB 0 ($-a_\alpha = -a_\beta$). We then negate the negation under the new encoding ($-a_\beta \mapsto a_\beta$), obtaining the original value in the new encoding.

As can easily be seen, the overhead we incur in addition to the cost of a Sign-Magnitude multiplication for multiplying two numbers of lengths n and m is basically that of 3 Two's Complement inversions, 3 Sign-Magnitude inversions (both of lengths n, m and $n + m$), and the cost of multiplying the boolean values representing whether the different cases are true or false. In total, the overhead costs (i.e., those incurred in addition to the costs for the Sign-Magnitude multiplication) are:

- $14(n + m) - 7$ binary additions
- $6(n + m) - 3$ binary multiplications
- a multiplicative depth of $\max\{n, m\} + 1 + n + m$

We present some concrete values for this overhead and runtimes in Fig. 3 along with the same values for Two's Complement multiplication and Sign-Magnitude multiplication. As can easily be seen, HE performs better than Two's Complement in all aspects for multiplying large numbers, but is (naturally) not quite as good as Sign-Magnitude. The runtimes are roughly as we would expect from these numbers, i.e., the new multiplication is faster than Two's Complement for large numbers, but naturally slower than Sign-Magnitude.

Thus, we have found a new way to improve efficiency for large bitlengths: do all operations in Two's Complement notation, but switch to Sign-Magnitude for multiplication. We shall see the benefits of this in our real-world application in Sect. 6.3, though we would like to note that there may be applications where Sign-Magnitude is favorable (when there are very few additions). However, since in Fully Homomorphic Encryption, multiplicative depth is often key (as mentioned in Sect. 3.1) and bootstrapping is the bottleneck, our new approach seems favorable for large parameters under this aspect as well.

5.2 Easy Comparison

Apart from numerical computations, many algorithms require a comparison of two numbers, which would usually require a rather expensive computation. However, we argue that in some use cases where one only has to compare a number to 0, like in the Perceptron, there is a much easier way. Instead of computing a costly circuit for comparison, it suffices to take the most significant bit of the number, which will be 0 if the number is greater than zero and 1 if it is less. For Two's Complement, it will be 0 also when the number equals 0, but in Sign-Magnitude it can be either 0 or 1 when using this method, as there are two encodings of 0 here. Thus, if the sum is exactly 0, the resulting bit is wrong for Two's Complement and can be either case for Sign-Magnitude. We observe, however, that when initializing the weights w_1, \dots, w_l with random rational numbers, a weighted sum $w_1x_1 + \dots + w_lx_l$ is highly unlikely to be 0. Thus, in this case there should be no change whether the condition for an operation is $w_1x_1 + \dots + w_lx_l > 0$ or $w_1x_1 + \dots + w_lx_l \geq 0$ and the easy comparison

should return the correct result with overwhelming probability. If the weights are initialized with 0 (as could be chosen in the Perceptron) or integers in the more general case, a more involved formula should be used.

5.3 Improved Multiplication

As the reader may have noticed, the sign extension in Two's Complement introduces costly redundancy, which can be avoided by carefully copying values to appropriate locations instead of computing them from scratch every time. Of course, as Sign-Magnitude multiplication works without sign extension, this improvement only applies to Two's Complement. However, the following further improvements hold for both encodings:

Having computed the matrix whose rows we want to sum up, we can apply a $\log(n+m)$ -depth circuit for adding the $n+m$ rows. It is noteworthy that we can save computation power by modifying the addition operation: As can easily be seen, we are always adding rows of different lengths. While the naive approach of padding the right-hand side of the shorter number with 0's and applying normal addition would also work, we can save some effort by copying the excess bits of the longer number and then performing addition on the remaining shorter equal-length parts. Generally, when using this second approach, we only perform an addition of the length of the shorter input, which is an important factor in depth optimization.

In the simpler case where one value is known, i.e., multiplication by a constant, we do not need to do as much work: For simplicity, assume that the input b is known. We again first need to do sign extension for Two's Complement, but in the next step instead of having to compute $n \cdot m$ terms $a_i \cdot b_j$ as before, we can just copy the string a for every bit that is 1 in b , shifting to the left with each bit. This way, we save $n \cdot m$ multiplications from the generation of the matrix and reduce the depth by one. Also, note that we now don't need to add as many rows, as we only write down those that correspond to the non-zero bits in b . Thus, we only need to do $\text{hm}(b)$ row additions, where $\text{hm}(b)$ is the hamming weight of b . Of course, the complexity and multiplicative depth now depend on the value of b and are the same as for regular multiplication in the worst case. However, on average we will only have to do half as many row additions.

5.4 Managing Length

By default, each addition and each multiplication increase the bitlength: Addition increases it by 1, whereas multiplication results in a bitlength that is the sum of the two input lengths. When performing several multiplications consecutively, this can easily lead to enormous bitlengths. However, in a scenario where the size of the values can be estimated, there is a way around this. One such scenario is machine learning, where the person working on the data is the person who has the algorithm for building the model and it is a reasonable assumption that some factors of the model are known, e.g. from experience. For example, in the data set we worked with [2], the value w_0 always took some value near

10000 no matter what subset of test subjects we chose. In such cases, the service provider who is doing the computations can put a bound on the lengths (i.e., he is certain that the weights will not be larger in absolute value than 2^q for some q). When this is the case, we can reduce the bitlength of the encrypted values to this size $q + 1$ by discarding the excess bits: In Two’s Complement, we can delete the most significant bits (which will all be 0 for a positive and 1 for a negative number) until we reach the desired length, whereas for Sign-Magnitude we discard the bits following the MSB (which will all be 0). More specifically, we actually integrated this into our multiplication routine, such that we not only save space, but also effort, as we only compute until we reach the bound in each step. This can be viewed as the inversion of the sign extension operation introduced in Sect. 3.3 and makes the entire algorithm significantly faster, as we have eliminated linear growth in the bitlength.

6 Applications

In this section, we demonstrate the performance increase on two concrete use cases.

6.1 Background and Motivation

Fully Homomorphic Encryption allows the computation of arbitrary functions on encrypted data while keeping the data hidden from the computing party. While FHE does not in principle offer to keep the function private (e.g., if the data and the function belong to the same party, who wishes to have the computation done by a different party with more computing power), it can hide the function that was applied in the following case: If the data belongs to one party and the function belongs to the computing party, then FHE schemes that are “circuit private” guarantee that a ciphertext divulges nothing about the function that was applied to it. Since circuit privacy is often a goal for FHE schemes, it makes sense to extend this requirement to the encoding choices to achieve privacy for the end result. This then means that the data owner learns nothing about the applied function except for what he can derive from the result, and the function owner learns nothing about the data. In this spirit, machine learning has often been cited as an application of Fully Homomorphic Encryption (see Sect. 2). Machine learning describes a field of research focused on extracting information from data, e.g. in the form of models. In this paper we consider the following scenario: Suppose Alice has a machine learning algorithm which takes data as input and returns a predictive model, and Bob has some data and would like either to obtain a model based on his data, or apply said model to further data (though he does not obtain the model in that case, e.g. allowing the service provider to bill him for each classification of his data). However, Alice does not want to reveal her algorithm for building the model to Bob, and Bob wishes to keep his data secret. With Fully Homomorphic Encryption, Bob could encrypt his (training) data and send it to Alice, who then performs her algorithm on

the encrypted data. The output is an encryption of the model, which Alice can apply to new encrypted data instances from Bob and Bob only receives the result of applying the model to his data (first case), or the whole model is sent to Bob (second case), in which case only Bob can decrypt the model. Thus, with an adequately secure Fully Homomorphic Encryption scheme, Alice has learned nothing about Bob's data and Bob has learned nothing about Alice's algorithm except what he can deduce from the result of the evaluation.

In the following, we consider two use cases, one for each of the above scenarios. For the first case, we take up a use case already presented in [22]: the Linear Means Classifier, where we assume that the model has already been built. Alice receives Bob's encrypted data, which she classifies by evaluating a polynomial of degree 2. This use case showcases our new Hybrid Encoding, which performs significantly better in this general case where the results are not bounded.

For the second case, we examine the Perceptron and show how to improve efficiency in evaluating it (i.e., obtaining the model), showcasing our results regarding choice of encoding and tweaks in multiplication. The Perceptron is an important fundamental algorithm in machine learning upon which many others are built, so being able to efficiently homomorphically evaluate it is mandatory before we can move on to more advanced machine learning algorithms.

The given runtimes are estimates for data encrypted with HElib [1], as runtimes are still very large: We measured the time for operations like addition and multiplication for different parameters and extrapolated the time it would take to compute the entire function. For example, given the function $f(x_1, x_2, x_3, x_4) = x_1 \cdot x_2 + x_3 \cdot x_4$ on inputs of length n , we would calculate the runtime as that of 2 multiplications of n -bit numbers plus one addition of numbers of lengths $2n$ (in the unbounded case). We confirmed our computations by actually running the Perceptron for lengths $n = 3$ and $n = 5$ for all three encodings to make sure that our computations reflect reality. However, we point out that these runtimes depend greatly on the characteristics of HElib: If one used a different encryption scheme that takes longer or shorter to perform bootstrapping, the results would vary greatly. However, our theoretical results are independent of the scheme that was used.

6.2 Linear Means Classifier

In this section, we examine the Linear Means Classifier to showcase the first use case, where the Service Provider retains the encrypted model and the user may send further encrypted data which is then classified by the encrypted model and only the encrypted result is returned to the user.

The Linear Means Classifier: Like [22], we consider the case where there are two classes, which are determined by the sign of the score function, which is a polynomial of degree 2. More concretely, the model consists of a vector $w = (w_1, \dots, w_l)$ and a constant c , and the data to be classified is a l -dimensional real-valued vector $x = (x_1, \dots, x_l)$. The score function is then computed as $\langle w, x \rangle + c = w_1x_1 + w_2x_2 + \dots + w_lx_l + c$, and the sign of the result determines

which class the data instance belongs to. As can easily be seen, this is closely related to the classification function of the Perceptron from the next section, where the focus is on determining w and c instead of computing the score function for given (encrypted and thus unknown) values for w and c as we do here.

Performance: Using the Linear Means Classifier, we examine the effects of using different encodings in the unbounded case (i.e., when the product of two n -bit numbers has length $2n$). To this end, we compute both the effort required in terms of bit operations and depth and the runtime of evaluating the score function for inputs of bitlength 30 for different numbers l of features. As explained above, we computed these runtimes from their components (i.e., the runtime for multiplying two 30-bit numbers without bounds, and the runtime for adding two 60-bit numbers) as the numbers are quite large. The results can be found in Fig. 2 in Sect. 1.1. As we can see, Two’s Complement is better than Sign-Magnitude, and using our new Hybrid Encoding significantly improves all aspects except depth, which is about halfway between the other two encodings. This did not matter in our case as we bootstrapped after every multiplication.

6.3 Homomorphically Evaluating the Perceptron

In this section, we examine the first use case where the Perceptron is evaluated to return an encrypted model.

The Perceptron: The Perceptron is an algorithm based on neural networks and basically works by computing a weighted sum of the input traits (usually rational numbers) for each subject and then classifying into one of two classes depending on whether this weighted sum is above a certain threshold or not. In the training phase, the weights are adjusted if the computed classification does not match the known classification of the training instance. After training, the model can be used to classify future inputs with no known classification. The model consists of the weights, and the threshold can either be predetermined or flexible (and thus part of the model being computed). We will work with the latter approach, which enables us to compare the inner product to 0.

Performance: We will now examine how the optimizations from Sect. 5 affect the Perceptron, as shown in Fig. 4. We can see that bounding the values makes a huge difference, especially since these values are only for the first round and would grow exponentially in further rounds. Sign-Magnitude is consistently the worst choice, and in the unbounded case, Hybrid Encoding is fastest (as already evident from Sect. 6.2). In the bounded case, however, Two’s Complement is fastest, and this makes sense: The fact that we have integrated the bounding into our multiplication procedure and stop computing in each line as soon as the bound is reached negates the sign extension that incurs the slowdown for multiplication in Two’s Complement encoding. This means that we expect bounded Two’s Complement multiplication to be almost as fast as Sign-Magnitude multiplication, which was confirmed by our experiments. Due to this, there is no efficiency gain through our new encoding in the bounded case, but the graph still

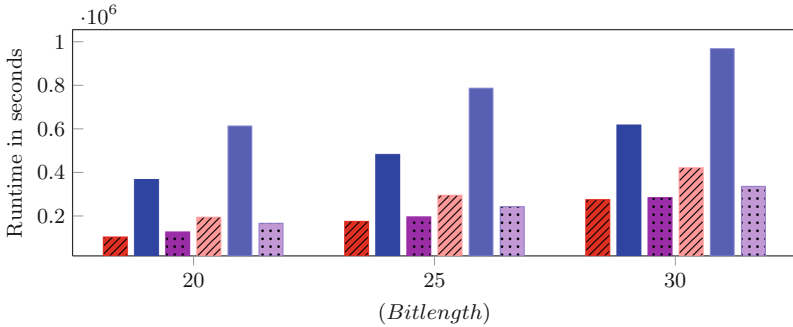


Fig. 4. Extrapolated runtimes for one subject for one round of the encrypted Perceptron for Two's Complement (• (lines) for bounded, • (lines) for unbounded values), Sign-Magnitude (• (solid) for bounded, • (solid) for unbounded values) and using our new Hybrid Encoding (• (dotted) for bounded, • (dotted) for unbounded values) (Color figure online).

illustrates the importance of choosing the right encoding, as Sign-Magnitude is significantly slower here due to its costly addition.

7 Conclusion and Future Work

In conclusion, we have presented a way of working with encrypted rational numbers, to our knowledge being the first to not restrict ourselves to unsigned integers. We have presented a new hybrid encoding technique that vastly improves efficiency for FHE on rational numbers both in theory and for real-world applications like the Linear Means Classifier, and other optimizations that improve efficiency for more complicated functions like the Perceptron. Since our results are independent of the scheme used, they hold with maximum generality and can thus be beneficial for anyone looking to evaluate a function homomorphically. For future research, we believe that this hybrid approach may be transferable to plaintext spaces other than $\{0, 1\}$, although the elementary operations will be considerably more involved. Further, we imagine that it could be beneficial to take a step back from established encodings and come up with a new one from scratch, which could be specially tailored to FHE computations.

References

1. HeLib Library: <https://github.com/shaih/HElib>
2. Pima Dataset: <https://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>
3. Armknecht, F., Boyd, C., Carr, C., Gjøsteen, K., Jäschke, A., Reuter, C.A., Strand, M.: A guide to fully homomorphic encryption. IACR Cryptology ePrint Archive (2015/1192)
4. Armknecht, F., Katzenbeisser, S., Peter, A.: Group homomorphic encryption: characterizations, impossibility results, and applications. DCC **67**(2), 209–232 (2013)

5. Armknecht, F., Strufe, T.: An efficient distributed privacy-preserving recommendation system. In: Med-Hoc-Net (2011)
6. Aslett, L.J.M., Esperança, P.M., Holmes, C.C.: Encrypted statistical machine learning: new privacy preserving methods. CoRR abs/1508.06845 (2015)
7. Bos, J.W., Lauter, K.E., Naehrig, M.: Private predictive analysis on encrypted medical data. *J. Biomed. Inform.* **50**, 234–243 (2014)
8. Bost, R., Popa, R.A., Tu, S., Goldwasser, S.: Machine learning classification over encrypted data. In: NDSS (2015)
9. Brakerski, Z., Gentry, C., Vaikuntanathan, V.: Fully homomorphic encryption without bootstrapping. *ECCC* **18**, 111 (2011)
10. Brakerski, Z., Vaikuntanathan, V.: Efficient fully homomorphic encryption from (standard) LWE. In: FOCS (2011)
11. Cheon, J.H., Kim, M., Lauter, K.: Homomorphic computation of edit distance. In: Brenner, M., Christin, N., Johnson, B., Rohloff, K. (eds.) FC 2015 Workshops. LNCS, vol. 8976, pp. 194–212. Springer, Heidelberg (2015)
12. Chung, H., Kim, M.: Encoding rational numbers for fhe-based applications. IACR Cryptology ePrint Archive (2016/344)
13. Coron, J.-S., Lepoint, T., Tibouchi, M.: Scale-invariant fully homomorphic encryption over the integers. In: Krawczyk, H. (ed.) PKC 2014. LNCS, vol. 8383, pp. 311–328. Springer, Heidelberg (2014)
14. Coron, J.-S., Naccache, D., Tibouchi, M.: Public key compression and modulus switching for fully homomorphic encryption over the integers. In: Pointcheval, D., Johansson, T. (eds.) EUROCRYPT 2012. LNCS, vol. 7237, pp. 446–464. Springer, Heidelberg (2012)
15. Costache, A., Smart, N.P., Vivek, S., Waller, A.: Fixed point arithmetic in SHE scheme. IACR Cryptology ePrint Archive (2016/250)
16. van Dijk, M., Gentry, C., Halevi, S., Vaikuntanathan, V.: Fully homomorphic encryption over the integers. In: Gilbert, H. (ed.) EUROCRYPT 2010. LNCS, vol. 6110, pp. 24–43. Springer, Heidelberg (2010)
17. Dowlin, N., Gilad-Bachrach, R., Laine, K., Lauter, K., Naehrig, M., Wernsing, J.: Manual for using homomorphic encryption for bioinformatics. Technical report MSR-TR-2015-87, Microsoft Research (2015)
18. Ducas, L., Micciancio, D.: FHEW: bootstrapping homomorphic encryption in less than a second. In: Oswald, E., Fischlin, M. (eds.) EUROCRYPT 2015. LNCS, vol. 9056, pp. 617–640. Springer, Heidelberg (2015)
19. Gentry, C.: A fully homomorphic encryption scheme. Ph.D. thesis, Stanford University (2009)
20. Gentry, C., Halevi, S., Smart, N.P.: Homomorphic evaluation of the AES circuit. In: Canetti, R., Safavi-Naini, R. (eds.) CRYPTO 2012. LNCS, vol. 7417, pp. 850–867. Springer, Heidelberg (2012)
21. Gentry, C., Sahai, A., Waters, B.: Homomorphic encryption from learning with errors: conceptually-simpler, asymptotically-faster, attribute-based. In: Canetti, R., Garay, J.A. (eds.) CRYPTO 2013, Part I. LNCS, vol. 8042, pp. 75–92. Springer, Heidelberg (2013)
22. Graepel, T., Lauter, K., Naehrig, M.: ML confidential: machine learning on encrypted data. In: Kwon, T., Lee, M.-K., Kwon, D. (eds.) ICISC 2012. LNCS, vol. 7839, pp. 1–21. Springer, Heidelberg (2013)
23. Henecka, W., Kögl, S., Sadeghi, A., Schneider, T., Wehrenberg, I.: TASTY: tool for automating secure two-party computations. In: CCS (2010)
24. Naehrig, M., Lauter, K.E., Vaikuntanathan, V.: Can homomorphic encryption be practical? In: CCSW (2011)

25. Sadeghi, A.-R., Schneider, T.: Generalized universal circuits for secure evaluation of private functions with application to data classification. In: Lee, P.J., Cheon, J.H. (eds.) ICISC 2008. LNCS, vol. 5461, pp. 336–353. Springer, Heidelberg (2009)
26. Smart, N.P., Vercauteren, F.: Fully homomorphic encryption with relatively small key and ciphertext sizes. In: Nguyen, P.Q., Pointcheval, D. (eds.) PKC 2010. LNCS, vol. 6056, pp. 420–443. Springer, Heidelberg (2010)
27. Songhori, E.M., Hussain, S.U., Sadeghi, A., Schneider, T., Koushanfar, F.: Tinygarble: highly compressed and scalable sequential garbled circuits. In: SP (2015)
28. Wu, D.J., Feng, T., Naehrig, M., Lauter, K.E.: Privately evaluating decision trees and random forests. IACR Cryptology ePrint Archive (2015/386)