



Efficient Homomorphic Comparison Methods with Optimal Complexity

Jung Hee Cheon^{1,2}, Dongwoo Kim¹, and Duhyeong Kim¹(✉)

¹ Department of Mathematical Sciences, Seoul National University,
Seoul, Republic of Korea

{jhcheon, dwkim606, doodoo1204}@snu.ac.kr

² Crypto Lab Inc., Seoul, Republic of Korea

Abstract. Comparison of two numbers is one of the most frequently used operations, but it has been a challenging task to efficiently compute the comparison function in homomorphic encryption (HE) which basically supports addition and multiplication. Recently, Cheon et al. (Asiacrypt 2019) introduced a new approximate representation of the comparison function with a rational function, and showed that this rational function can be evaluated by an iterative algorithm. Due to this iterative feature, their method achieves a logarithmic computational complexity compared to previous polynomial approximation methods; however, the computational complexity is still not optimal, and the algorithm is quite slow for large-bit inputs in HE implementation.

In this work, we propose new comparison methods with *optimal* asymptotic complexity based on *composite polynomial* approximation. The main idea is to systematically design a constant-degree polynomial f by identifying the *core properties* to make a composite polynomial $f \circ f \circ \dots \circ f$ get close to the sign function (equivalent to the comparison function) as the number of compositions increases. We additionally introduce an acceleration method applying a mixed polynomial composition $f \circ \dots \circ f \circ g \circ \dots \circ g$ for some other polynomial g with different properties instead of $f \circ f \circ \dots \circ f$. Utilizing the devised polynomials f and g , our new comparison algorithms only require $\Theta(\log(1/\epsilon)) + \Theta(\log \alpha)$ computational complexity to obtain an approximate comparison result of $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$ within $2^{-\alpha}$ error.

The asymptotic optimality results in substantial performance enhancement: our comparison algorithm on 16-bit encrypted integers for $\alpha = 16$ takes 1.22 ms in amortized running time based on an approximate HE scheme HEAAN, which is 18 times faster than the previous work.

1 Introduction

Homomorphic Encryption (HE) is a primitive of cryptographic computing, which allows computations over encrypted data without any decryption process. With HE, clients who sent encrypted data to an untrusted server are guaranteed data privacy, and the server can perform any operations over the encrypted data. In recent years, HE has gained worldwide interest from various fields related to data

privacy issues including genomics [37–39] and finances [3, 31]. In particular, HE is emerging as one of the key tools to protect data privacy in machine learning tasks, which now became a necessary consideration due to public awareness of data breaches and privacy violation.

The comparison function $\text{comp}(a, b)$, which outputs 1 if $a > b$, 0 if $a < b$ and $1/2$ if $a = b$, is one of the most prevalent operations along with addition and multiplication in various real-world applications. For example, many of the machine learning algorithms such as cluster analysis [17, 33], gradient boosting [25, 26], and support-vector machine [19, 40] require a number of comparison operations. Therefore, it is indispensable to find an efficient method to compute the comparison function in an encrypted state for HE applications.

Since HE schemes [7, 11, 24] basically support homomorphic addition and multiplication, to compute non-polynomial operations including the comparison function in an encrypted state, we need to exploit polynomial approximations on them. The usual polynomial approximation methods such as minimax find approximate polynomials with minimal degree on a target function for given a certain error bound. However, the computational complexity to evaluate these polynomials is so large that it is quite inefficient to obtain approximate results with high-precision by these methods. Recently, to resolve this problem, Cheon et al. [12] introduced a new identity $\text{comp}(a, b) = \lim_{k \rightarrow \infty} a^k / (a^k + b^k)$, and showed that the identity can be computed by an iterative algorithm. Due to this iterative feature, their algorithm achieves a logarithmic computational complexity compared to usual polynomial approximation methods. However, the algorithm only achieves quasi-optimal computational complexity, and it is quite slow in HE implementation; more than 20 min is required to compute a single homomorphic comparison of 16-bit integers.

In this work, we propose new comparison methods using composite polynomial approximation on the sign function, which is equivalent to the comparison function. Starting from the analysis on the behavior of a composite polynomial $f^{(d)} := f \circ f \circ \dots \circ f$, we identify the *core properties* of f that make $f^{(d)}$ get close to the sign function as d increases. We additionally introduce a novel acceleration method by applying a *mixed composition* of f and some other polynomial g with different properties instead of a simple composition of f . Applying these systematically devised polynomials f and g , we construct new comparison algorithms which firstly achieve the *optimal* computational complexity *among all polynomial evaluations* to obtain an approximate value of the comparison result within a certain error bound.

Our composite polynomial methods can be directly applied to evaluate piecewise polynomials with two sub-polynomials including the absolute function: For example, the function p such that $p(x) = p_1(x)$ if $x \in [0, 1]$ and $p(x) = p_2(x)$ if $x \in [-1, 0]$ for polynomials p_1 and p_2 can be represented by $p_1(x) \cdot (1 + \text{sgn}(x))/2 + p_2(x) \cdot (1 - \text{sgn}(x))/2$. Furthermore, our method is potentially applicable to more general piecewise polynomials including step functions (see Remark 1).

1.1 Our Idea and Technical Overview

Our key idea to identify several core properties of the basic function f essentially comes from a new interpretation of the previous work [12]. To be precise, [12] exploits the following identity to construct a comparison algorithm:

$$\lim_{k \rightarrow \infty} \frac{a^k}{a^k + b^k} = \begin{cases} 1 & \text{if } a > b \\ 1/2 & \text{if } a = b \\ 0 & \text{if } a < b \end{cases} = \text{comp}(a, b)$$

for positive numbers $a, b \in [1/2, 3/2]$. Since very large exponent $k = 2^d$ is required to obtain a comparison result within small error, they suggest to iteratively compute $a \leftarrow a^2/(a^2 + b^2)$ and $b \leftarrow b^2/(a^2 + b^2)$ with an initial step $a \leftarrow a/(a + b)$ and $b \leftarrow b/(a + b)$, which results in $a^{2^d}/(a^{2^d} + b^{2^d}) \simeq \text{comp}(a, b)$ after d iterations. The inverse operation $1/(a^2 + b^2)$ in each iteration is computed by Goldschmidt’s division algorithm [30].

The computational inefficiency of the comparison algorithm in [12] mainly comes from that inverse operation which should be done at least d times. Then, the natural question would be

“How can we construct an efficient comparison algorithm without inverse operation?”

To do this, we analyze the comparison algorithm in [12] with a new perspective. Let $f_0(x) = x^2/(x^2 + (1 - x)^2)$, then each iteration $a \leftarrow a^2/(a^2 + b^2)$ and $b \leftarrow b^2/(a^2 + b^2)$ can be interpreted as an evaluation of $f_0(a)$ and $f_0(b) = 1 - f_0(a)$ for $0 \leq a, b \leq 1$, respectively. Indeed, the d iterations correspond to the d -time composition of the basic function f_0 denoted by $f_0^{(d)} := f_0 \circ f_0 \circ \dots \circ f_0$, and the comparison algorithm can be interpreted as approximating $(\text{sgn}(2x - 1) + 1)/2$ by a composite polynomial $f_0^{(d)}$ (Fig. 1).

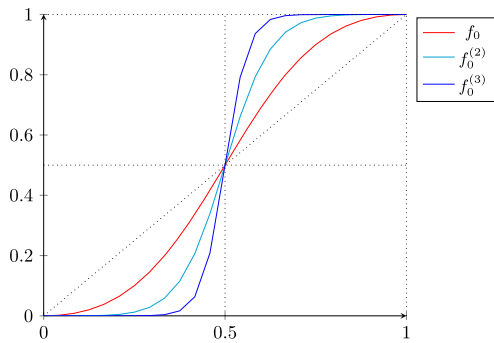


Fig. 1. Illustration of $f_0^{(d)}$ for $d = 1, 2, 3$

Our key observation on the basic function f_0 is that we actually do not need the exact formula of $f_0(x) = x^2/(x^2 + (1 - x)^2)$. Instead, it suffices to use

other polynomials with *similar shape* to f_0 : convex in $[0, 0.5]$, concave in $[0.5, 1]$, symmetric to the point $(0.5, 0.5)$, and have a value 1 at $x = 1$. For example, the composition $h_1^{(d)}$ of our devised polynomial $h_1(x) = -2x^3 + 3x^2$, which has similar shape to f_0 , gets close to $(\text{sgn}(2x - 1) + 1)/2$ as d increases. As a result, we can approximate the comparison function by a composite polynomial $f^{(d)}$ for some constant-degree polynomial f with several *core properties*, and identifying these core properties is the most important step in our algorithm construction.

Core Properties of f . Since the sign function is equivalent to the comparison function, via $\text{sgn}(x) = 2 \cdot \text{comp}(x, 0) - 1$ and $\text{comp}(a, b) = (\text{sgn}(a - b) + 1)/2$, it is enough to find a polynomial f such that $f^{(d)}(x)$ gets close to $\text{sgn}(x)$ over $[-1, 1]$ for some proper d . The core properties of f are as following:

Prop I. $f(-x) = -f(x)$

Prop II. $f(1) = 1, f(-1) = -1$

Prop III. $f'(x) = c(1 - x)^n(1 + x)^n$ for some constant $c > 0$

The first property is necessary from the origin symmetry of the sign function, and the second property is required to achieve $\lim_{d \rightarrow \infty} f^{(d)}(x) = 1$ for $0 < x \leq 1$. The last property makes f to be concave in $[0, 1]$ and convex in $[-1, 0]$, and the multiplicity n of ± 1 in $f'(x)$ accelerates the convergence of $f^{(d)}$ to the sign function. Interestingly, for each $n \geq 1$, a polynomial f_n satisfying above three properties is uniquely determined as

$$f_n(x) = \sum_{i=0}^n \frac{1}{4^i} \cdot \binom{2i}{i} \cdot x(1 - x^2)^i.$$

Since $\text{sgn}(x)$ is a discontinuous function at $x = 0$, the closeness of a polynomial $f(x)$ to $\text{sgn}(x)$ should be considered carefully. Namely, we do not consider a small neighborhood $(-\epsilon, \epsilon)$ of zero when measuring the difference between $f(x)$ and $\text{sgn}(x)$ (if not, the infinite norm is always ≥ 1). In Sect. 3.2, we prove that the infinite norm of $f_n^{(d)}(x) - \text{sgn}(x)$ over $[-1, -\epsilon] \cup [\epsilon, 1]$ is smaller than $2^{-\alpha}$ if $d \geq d_n$ for some $d_n > 0$. Then, $(f_n^{(d_n)}(a - b) + 1)/2$ outputs an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$.

Acceleration Method. Along with $\{f_n\}_{n \geq 1}$, we provide another family of odd polynomials $\{g_n\}_{n \geq 1}$ which reduces the required number of polynomial compositions d_n . At a high-level, we can interpret d_n as $d_n := d_\epsilon + d_\alpha$ where each of the terms d_ϵ and d_α has distinct aim as following: The first term d_ϵ is a required number of compositions to map the interval $[\epsilon, 1]$ into the interval $[1 - \tau, 1]$ for some fixed constant $0 < \tau < 1$ (typically, $\tau = 1/4$), and the second term d_α is a required number of compositions to map $[1 - \tau, 1]$ into $[1 - 2^{-\alpha}, 1]$, i.e.,

$$\begin{aligned} f_n^{(d_\epsilon)}([\epsilon, 1]) &\subseteq [1 - \tau, 1], \\ f_n^{(d_\alpha)}([1 - \tau, 1]) &\subseteq [1 - 2^{-\alpha}, 1]. \end{aligned}$$

In this perspective, our idea is to reduce d_ϵ by substituting $f_n^{(d_\epsilon + d_\alpha)}$ with $f_n^{(d_\alpha)} \circ g_n^{(d_\epsilon)}$ for some other $(2n + 1)$ -degree polynomial g_n with *weaker properties* than

the core properties of f_n . Since the first d_ϵ compositions only need to map $[\epsilon, 1]$ into $[1 - \tau, 1]$, Prop II & III are *unnecessary* in this part. Instead, the following property along with Prop I is required:

Prop IV. $\exists 0 < \delta < 1$ s.t. $x < g_n(x) \leq 1$ for $x \in (0, \delta]$ and $g_n([\delta, 1]) \subseteq [1 - \tau, 1]$

For g_n satisfying Prop I & IV, the composition $g_n^{(d)}$ does not get close to the sign function as d increases; however, we can guarantee that $g_n^{(d_\epsilon)}([\epsilon, 1]) \subseteq [1 - \tau, 1]$ for some $d_\epsilon > 0$ which is exactly the aim of first d_ϵ compositions. With some heuristic properties on g_n obtained by Algorithm 2, the required number of the first-part compositions d_ϵ is reduced by nearly half (see Sect. 3.5).

1.2 Our Results

New Comparison Methods with Optimal Complexity. We first propose a family of polynomials $\{f_n\}_{n \geq 1}$ whose composition $f_n^{(d)}$ gets close to the sign function (in terms of (α, ϵ) -closeness) as d increases. Based on the approximation

$$\frac{f_n^{(d)}(a - b) + 1}{2} \simeq \frac{\text{sgn}(a - b) + 1}{2} = \text{comp}(a, b),$$

we construct a new comparison algorithm $\text{NewComp}(a, b; n, d)$ which achieves *optimal asymptotic complexity* among the polynomial evaluations obtaining an approximate value of comparison within a certain level of error. The following theorem is the first main result of our work:

Theorem 1. *If $d \geq \frac{2+o(1)}{\log n} \cdot \log(1/\epsilon) + \frac{1}{\log n} \cdot \log \alpha + O(1)$, the comparison algorithm $\text{NewComp}(a, b; n, d)$ outputs an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$.*

The theorem implies that one can obtain an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$ with $\Theta(\log(1/\epsilon)) + \Theta(\log \alpha) + O(1)$ complexity and depth with NewComp .

We also provide another family of polynomials $\{g_n\}_{n \geq 1}$, which enables to reduce the number of polynomial compositions by substituting $f_n^{(d)}$ with $f_n^{(d_f)} \circ g_n^{(d_g)}$. From the mixed polynomial composition, we construct another comparison algorithm NewCompG with the following result:

Theorem 2 (Heuristic). *If $d_g \geq \frac{1+o(1)}{\log n} \cdot \log(1/\epsilon) + O(1)$ and $d_f \geq \frac{1}{\log n} \cdot \log \alpha + O(1)$, the comparison algorithm $\text{NewCompG}(a, b; n, d_f, d_g)$ outputs an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$.*

Since g_n and f_n have the same degree, the total depth and computational complexity of NewCompG are strictly smaller than those of NewComp .

The variety on choosing n in our comparison algorithms provides flexibility in complexity-depth tradeoff. For instance, one can choose $n = 4$ to achieve the minimal computational complexity (see Sect. 3.4). On the other hand, if one

wants to obtain comparison results with larger complexity but smaller depth, one can choose n larger than 4. Assuming some heuristic properties of g_n , the total depth of $\text{NewCompG}(\cdot, \cdot; n, d_f, d_g)$ gets close to the theoretical minimal depth as n increases (see Sect. 3.5).

Improved Performance. For two 8-bit integers which are encrypted by an approximate HE scheme HEAAN [11], the comparison algorithm NewComp (for $\epsilon = 2^{-8}$ and $\alpha = 8$) takes 0.9 ms in amortized running time, and the performance is twice accelerated by applying the other comparison algorithm NewCompG . The implementation result on NewCompG is about 8 times faster than that on the comparison algorithm of the previous work [12] based on HEAAN. Note that this performance gap grows up as the bit-length of input integers increases: For two encrypted 20-bit integers, our algorithm NewCompG is about 30 times faster than the previous work.

Application to Max. Since the max function is expressed by the sign function as $\max(a, b) = \frac{a+b}{2} + \frac{a-b}{2} \cdot \text{sgn}(a-b)$, we can directly obtain max algorithms from the family of polynomials $\{f_n\}_{n \geq 1}$ (and hence $\{g_n\}_{n \geq 1}$). Our max algorithms NewMax and NewMaxG outperform the max algorithm in the previous work [12] in terms of both computational complexity and depth. To be precise, the max algorithm in [12] requires $4\alpha + O(1)$ depth and $6\alpha + O(1)$ complexity to obtain an approximate value of min/max of two numbers in $[0, 1]$ within $2^{-\alpha}$ error. In our case, the max algorithm NewMax applying f_4 only require $3.08\alpha + O(1)$ depth and complexity, and it can be even reduced to $1.54\alpha + 1.72 \log \alpha + O(1)$ by using the other max algorithm NewMaxG . In practice, for encrypted 20-bit integers our NewMaxG algorithm is 4.5 times faster than the max algorithm in [12].

Moreover, our max algorithms fundamentally solve a potential problem of the max algorithm in [12] when inputs are encrypted by HEAAN. When two input numbers are too close so that the difference is even smaller than approximate errors of HEAAN, then the max algorithm in [12] may output a totally wrong result; in contrast, our max algorithms works well for any inputs from $[0, 1]$.

1.3 Related Works

Numerical Analysis on the Sign Function. In the literature of numerical analysis, to the best of our knowledge, there exist two main approaches on the polynomial approximation of the sign function. One is to naively apply general polynomial approximation methods (Taylor, least squares, minimax, etc.), and the other is to apply Newton's root-finding algorithm on a function which has ± 1 as roots.

General polynomial approximation methods provide an approximate polynomial with *minimal degree* under a certain upper bound of the approximate error. However, the evaluation of such approximate polynomial requires at least $\Theta(\sqrt{\text{degree}})$ multiplications, which yields *super-large computational complexity* when we aim to obtain a high-precision approximation. For example, when we want to obtain an approximate polynomial of the sign function with α -bit precision over $[-1, -2^{-\alpha}] \cup [2^{-\alpha}, 1]$ via general polynomial approximation methods,

the required computational complexity is at least $\Theta(\sqrt{\alpha} \cdot 2^{\alpha/2})$ which is exponential to α (see Sect. 2.2 for more details). There have been recent works [8, 32] applying Chebyshev polynomial approximation (on the sine function) instead of the minimax polynomial approximation for better efficiency. However, the Chebyshev polynomial approximation method still requires exponential computational complexity with respect to α when it is applied to the sign function.

Newton's root-finding algorithm outputs an approximate value of roots of a function $r(x)$ by iteratively computing $x_{n+1} = x_n - \frac{r(x_n)}{r'(x_n)}$ for an initial point x_0 . That is, an iterative computation of the function $f(x) = x - \frac{r(x)}{r'(x)}$ gives an approximate value to one of the roots of r . The most simple choice of r to compute the sign function is $r(x) = 1 - x^2$ which derives $f(x) = \frac{1}{2} \cdot (x + \frac{1}{x})$ so-called Newton's method [34, 36]. There have also been several attempts to improve the convergence rate of this iterative method to the sign function by changing f to $f(x) = \frac{3x+x^3}{1+3x^2}$ (Halley's method [42]), $f(x) = \frac{5x+10x^3+x^5}{1+10x^2+5x^4}$ [18], and $f(x) = \frac{10x+98x^3+126x^5+22x^7}{1+42x^2+140x^4+70x^6+3x^8}$ [46].¹ However, all these methods commonly require the inverse operation, and additional polynomial approximation on inverse is required to apply these methods in HE as the previous work [12]. Consequently, these methods are much less efficient than our methods for the evaluation of the sign function in HE due to a number of expensive inverse operations.

There has been proposed another choice of r that makes f a polynomial as in this paper, so-called Newton-Schulz method [34, 36]. When we take $r(x) = 1 - 1/x^2$, the function f is expressed as $f(x) = \frac{x}{2} \cdot (3 - x^2)$ and we can obtain an approximate value of the sign function by the iterative computation of f . Interestingly, this function is one of our devised polynomials f_1 . However, we note that the design rationale of our methods, setting core properties of f that makes $f^{(d)}$ get close to the sign function as d increases, is totally different from that of the Newton's root-finding method. With Newton's method it is not clear at all how to generalize f_1 to f_n for $n > 1$ or how to obtain the intuition for devising other polynomials $\{g_n\}_{n \geq 1}$ for convergence acceleration. Our methods applying $\{f_n\}_{n > 1}$ and $\{g_n\}_{n \geq 1}$ achieve much less computational complexity and depth than the previous numerical method (see Sect. 3.4 and Sect. 3.5).

HE-Based Comparison Methods. There have been several works on comparison algorithms for HE schemes [7, 11, 24] basically supporting addition and multiplication. The most recent work was proposed by Cheon et al. [12] which exploits the identity $\text{comp}(a, b) = \lim_{k \rightarrow \infty} \frac{a^k}{a^k + b^k}$ for $a, b > 0$ with an iterative inverse algorithm. Their comparison algorithm requires $\Theta(\alpha \log \alpha)$ complexity, which is quasi-optimal, to obtain an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for $a, b \in [1/2, 3/2]$ satisfying $\max(a, b)/\min(a, b) \geq 1 + 2^{-\alpha}$.

¹ In fact, this line of work in numerical analysis aims to compute the matrix sign function [36] which is a more general object than the sign function in our context. An inverse operation is not much more costly than a multiplication in their (asymptotic) cost analysis and experiments, which is a crucial difference from HE which requires an additional costly polynomial approximation for inverse [12].

There have been several approaches to approximate the sign function by polynomials to obtain a comparison algorithm. In 2018, Boura et al. [5] proposed an analytic method to compute the sign function by approximating it via Fourier series over a target interval which has an advantage on numerical stability. In this method, one should additionally consider the error induced by the polynomial approximation on e^{ix} . Another approach is to approximate the sign function by $\tanh(kx) = \frac{e^{kx} - e^{-kx}}{e^{kx} + e^{-kx}}$ for sufficiently large $k > 0$ [14]. In order to efficiently compute $\tanh(kx)$, they repeatedly apply the double-angle formula $\tanh(2x) = \frac{2 \tanh(x)}{1 + \tanh^2(x)} \approx \frac{2x}{1+x^2}$ where the inverse operation is substituted by a low-degree minimax approximate polynomial. This procedure can be interpreted as a composition of polynomial f which is the low-degree minimax approximation polynomial of $\frac{2x}{1+x^2}$. However, their method does not catch core properties of the basic polynomial f (e.g., $f(1) = 1$), so the error between $f^{(d)}$ and $\text{sgn}(x)$ cannot be reduced below a certain bound even if we increase d to ∞ . As an independent work, Bajard et al. [4] recently proposed a new approach to approximately compute the sign function by applying the Newton’s root-finding method on the function $r(x) = 1 - 1/x^2$, which corresponds to one of our devised polynomials f_1 .

When each bit of message is encrypted separately [13, 16, 20], one can perform a comparison operation of two α -bit integers with $O(\log \alpha)$ depth and $O(\alpha)$ complexity. The bit-by-bit encryption method was recently generalized to encrypt an integer a after decomposing it as $a = \sum a_i b^i$ for a power of small prime $b = p^r$ [47]. However, since these encryption methods are quite inefficient for addition and multiplication, they are not desirable when comparison operations are mixed with a number of polynomials such as cluster analysis and gradient tree boosting.

2 Preliminaries

2.1 Notations

All logarithms are of base 2 unless otherwise indicated, and e denotes the Euler’s constant. \mathbb{Z} , \mathbb{R} and \mathbb{C} denote the integer ring, the real number field and complex number field, respectively. For a finite set X , we denote the uniform distribution over X by $U(X)$. For a real-valued function f defined over \mathbb{R} and a compact set $I \subset \mathbb{R}$, we denote the infinity norm of f over the domain I by $\|f\|_{\infty, I} := \max_{x \in I} |f(x)|$. The d -times composition of f is denoted by $f^{(d)} := f \circ f \circ \dots \circ f$. We denote the sign function and the comparison function by

$$\text{sgn}(x) := \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x = 0 \\ -1 & \text{if } x < 0 \end{cases}, \quad \text{comp}(a, b) := \begin{cases} 1 & \text{if } a > b \\ 1/2 & \text{if } a = b \\ 0 & \text{if } a < b \end{cases},$$

which are in fact equivalent to each other by $\text{comp}(a, b) = (\text{sgn}(a - b) + 1)/2$.

For $\alpha > 0$ and $0 < \epsilon < 1$, we say a polynomial f is (α, ϵ) -close to $\text{sgn}(x)$ over $[-1, 1]$ if it satisfies

$$\|f(x) - \text{sgn}(x)\|_{\infty, [-1, -\epsilon] \cup [\epsilon, 1]} \leq 2^{-\alpha}.$$

For $a, b \in \mathbb{R}$, we denote the complexity $a \cdot \log(1/\epsilon) + b \cdot \log \alpha + O(1)$ by $L(a, b)$. The O notation in this paper regards to α and $1/\epsilon$. In the rest of this paper, we only consider the (non-scalar) multiplicative depth and (non-scalar) multiplicative computational complexity, i.e., we do not count the number of additions nor scalar multiplications in computational complexity.

2.2 Minimax Polynomial Approximation Method

In this paper, we measure the accuracy of polynomial approximation methods by the maximal error between the target function and an approximate polynomial over a predetermined domain. In this respect, the minimax approximation method provides the best approximate polynomials among general polynomial approximation methods. For a positive odd integer k , let us denote by $p_{k,\epsilon}$ the degree- k polynomial p which minimizes $\|\text{sgn}(x) - p(x)\|_{\infty, [-1, -\epsilon] \cup [\epsilon, 1]}$. For the sign function $\text{sgn}(x)$, there exists a tight lower bound on the approximation error:

$$\lim_{k \rightarrow \infty} \sqrt{\frac{k-1}{2}} \cdot \left(\frac{1+\epsilon}{1-\epsilon}\right)^{\frac{k-1}{2}} \cdot \|\text{sgn}(x) - p_{k,\epsilon}(x)\|_{\infty, [-1, -\epsilon] \cup [\epsilon, 1]} = \frac{1-\epsilon}{\sqrt{\pi\epsilon}}$$

for $0 < \epsilon < 1$, which was proved by Eremenko and Yuditskii [23]. More general works on minimax polynomial approximation of piecewise analytic function have been proposed [2, 44], but [23] provides more tight and accurate results on error analysis for the sign function.

Assume that k is large enough so that the left-hand side $\sqrt{\frac{k-1}{2}} \cdot \left(\frac{1+\epsilon}{1-\epsilon}\right)^{\frac{k-1}{2}} \cdot \|\text{sgn}(x) - p_{k,\epsilon}(x)\|_{\infty, [-1, -\epsilon] \cup [\epsilon, 1]}$ is sufficiently close to the limit value. To bound the approximation error by $2^{-\alpha}$ for $\text{sgn}(x)$ over $[-1, -\epsilon] \cup [\epsilon, 1]$, the degree k should be chosen to satisfy

$$\sqrt{\frac{k-1}{2}} \cdot \left(\frac{1+\epsilon}{1-\epsilon}\right)^{\frac{k-1}{2}} \cdot \frac{\sqrt{\pi\epsilon}}{1-\epsilon} > 2^\alpha,$$

which implies that k should be at least $\Theta(\alpha/\epsilon)$ from the fact $\log\left(\frac{1+\epsilon}{1-\epsilon}\right) \approx \frac{\epsilon}{2}$ for small ϵ . Then, the evaluation of the polynomial $p_{k,\epsilon}$ requires at least $\log \alpha + \log(1/\epsilon) + O(1)$ depth and $\Theta\left(\sqrt{\alpha/\epsilon}\right)$ complexity applying the Paterson-Stockmeyer method [43] which is asymptotically optimal.

There exists a well-known theorem called the equioscillation theorem attributed to Chebychev, which specifies the *shape* of the minimax approximate polynomial $p_{k,\epsilon}$.

Lemma 1 (Equioscillation Theorem for sign function [23]). *Let $\text{sgn}(x)$ be the sign function (Sect. 2.1). For $k \geq 1$ and $0 < \epsilon < 1$, an odd polynomial $p_{k,\epsilon}$ of degree $(2k + 1)$ minimizes the infinity norm $\|\text{sgn} - p_{k,\epsilon}\|_{\infty, [-1, -\epsilon] \cup [\epsilon, 1]}$ if and only if there are $k + 2$ points $\epsilon = x_0 < x_1 < \dots < x_{k+1} = 1$ such that $\text{sgn}(x_i) - p_{k,\epsilon}(x_i) = (-1)^i \|\text{sgn} - p_{k,\epsilon}\|_{\infty}$. Here, x_1, x_2, \dots, x_k are critical points.*

Note that the if-and-only-if statement of the above lemma also implies the *uniqueness* of the minimax polynomial approximation of $\text{sgn}(x)$ on $[-1, -\epsilon] \cup [\epsilon, 1]$ for given ϵ and degree $2k + 1$. In the rest of paper, we will use the fact that $p_{k,\epsilon}$ is concave and increasing in the interval $[0, x_0]$ (in fact it holds for $[0, x_1]$).

2.3 Homomorphic Encryption

HE is a cryptographic primitive which allows arithmetic operations including addition and multiplication over encrypted data without decryption process. HE is regarded as a promising solution which prevents leakage of private information during analyses on sensitive data (e.g., genomic data, financial data). A number of HE schemes [6, 7, 11, 15, 22, 24, 28] have been suggested following Gentry's blueprint [27], and achieving successes in various applications [5, 9, 29, 37].

In this paper, we mainly focus on word-wise HE schemes, i.e., the HE schemes whose basic operations are addition and multiplication of encrypted message vectors over $\mathbb{Z}/p\mathbb{Z}$ for $p \geq 2$ [7, 24, 28] or the complex number field \mathbb{C} [11]. An HE scheme consists of the following algorithms:

- KeyGen(params). For parameters `params` determined by a level parameter L and a security parameter λ , output a public key `pk`, a secret key `sk`, and an evaluation key `evk`.
- Enc_{pk}(`m`). For a message m , output the ciphertext `ct` of m .
- Dec_{sk}(`ct`). For a ciphertext `ct` of m , output the message m .
- Add_{evk}(`ct`₁, `ct`₂). For ciphertexts `ct`₁ and `ct`₂ of m_1 and m_2 , output the ciphertext `ct`_{add} of $m_1 + m_2$.
- Mult_{evk}(`ct`₁, `ct`₂). For ciphertexts `ct`₁ and `ct`₂ of m_1 and m_2 , output the ciphertext `ct`_{mult} of $m_1 \cdot m_2$.

Though any arithmetic circuit can be computed by HE theoretically, the number of multiplications and multiplicative depth of the circuit are major factors affecting the practical performance and feasibility in real-world applications.

3 Our New Comparison Method

Since the comparison function and the sign function are equivalent, it suffices to find a nice approximate polynomial (with one variable) of the sign function instead of the comparison function (with two variables). In this section, we will introduce new polynomial approximation methods for the sign function which we call *composite polynomial approximation*, and analyze their computational efficiency. As in [12], we assume that the input numbers are contained in the bounded interval $[0, 1]$, since $x \in [c_1, c_2]$ for known constants $c_1 < c_2$ can be scaled down into $[0, 1]$ via mapping $x \mapsto (x - c_1)/(c_2 - c_1)$. Therefore, the domain of $\text{sgn}(x)$ we consider in this paper is $[-1, 1]$.

3.1 Composite Polynomial Approximation of Sign Function

As described in [12], approximating a non-polynomial function by composite polynomials has an advantage in computational complexity: A composite function F of a constant-degree polynomial f , i.e., $F := f \circ f \circ \dots \circ f$, can be computed within $O(\log(\deg F))$ complexity, while the evaluation of an arbitrary polynomial G requires at least $\Theta(\sqrt{\deg G})$ [43]. However, even if this methodology achieves a log-degree computational complexity, it would be meaningless if the total degree of F is extremely large (e.g., $\deg F = 2^{\deg G}$). Therefore, it is very important to *well-design* a constant polynomial f so that it requires small d to make $f^{(d)}$ sufficiently close to $\text{sgn}(x)$ over $[-1, 1]$. Since $\text{sgn}(x)$ is discontinuous at $x = 0$, we are not able to obtain a nice polynomial approximation of $\text{sgn}(x)$ over $(-\epsilon, \epsilon)$ for any $0 < \epsilon < 1$. As a result, we set our goal to find f whose composition $f^{(d)}$ is (α, ϵ) -close to the sign function for $\alpha > 0$ and $0 < \epsilon < 1$ with small d .

The key observation for designing such polynomial f is as follows: For $x_0 \in [-1, 1]$, let x_i be the i -time composition value $f^{(i)}(x_0)$. Then, the behavior of x_i 's can be easily estimated with the graph of f . For example, given x_0 on the x -coordinate, x_1 can be identified by the x -coordinate of the intersection point of the graph $y = x$ and the horizontal line $y = f(x_0)$. Note that we can iteratively estimate x_{i+1} with the previous point x_i (see Fig. 2).

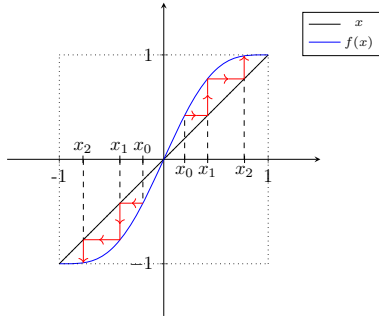


Fig. 2. Behavior of $x_i = f^{(i)}(x_0)$ for $f(x) = -\frac{5}{16}x^7 + \frac{21}{16}x^5 - \frac{35}{16}x^3 + \frac{35}{16}x$

In this perspective, the basic polynomial f should be constructed so that x_i gets close to 1 if $x_0 \in (0, 1]$ and -1 if $x_0 \in [-1, 0)$ as i increases. We can formally identify three properties of f as follows: Firstly, since the sign function is an odd function, we also set f to be an odd function. Secondly, we set $f(1) = 1$ and $f(-1) = -1$ to make $f^{(d)}(x)$ point-wise converge to $\text{sgn}(x)$ whose value is ± 1 for $x \neq 0$. More precisely, if $f^{(d)}(x)$ for some $x \in [-1, 1]$ converges to y as d increases, it must hold that $f(y) = f(\lim_{d \rightarrow \infty} f^{(d)}(x)) = \lim_{d \rightarrow \infty} f^{(d)}(x) = y$. Lastly, f should be considered as a *better* polynomial if it is *more concave* over $[0, 1]$ (hence *more convex* over $[-1, 0]$), which will accelerate the convergence of $f^{(d)}$ to the sign function. In order to increase convexity, we set the derivative

function f' of f to have maximal multiple roots at 1 and -1 . These properties are summarized as following.

Core Properties of f :

- Prop I. $f(-x) = -f(x)$ (Origin Symmetry)
- Prop II. $f(1) = 1, f(-1) = -1$ (Convergence to ± 1)
- Prop III. $f'(x) = c(1-x)^n(1+x)^n$ for some $c > 0$ (Fast convergence)

For a fixed $n \geq 1$, a polynomial f of the degree $(2n + 1)$ satisfying those three properties is uniquely determined, and we denote this polynomial by f_n (and the uniquely determined constant c by c_n): From Prop I and III, we get $f_n(x) = c_n \int_0^x (1-t^2)^n dt$, and the constant c_n is determined by Prop II. By applying the following identity

$$\int_0^x \cos^m t \, dt = \frac{1}{m} \cdot \cos^{m-1} x \cdot \sin x + \frac{m-1}{m} \cdot \int_0^x \cos^{m-2} t \, dt$$

which holds for any $m \geq 1$, we obtain

$$f_n(x) = \sum_{i=0}^n \frac{1}{4^i} \cdot \binom{2i}{i} \cdot x(1-x^2)^i.$$

See Appendix A for more details. Hence, we can easily compute f_n as following:

- $f_1(x) = -\frac{1}{2}x^3 + \frac{3}{2}x$
- $f_2(x) = \frac{3}{8}x^5 - \frac{10}{8}x^3 + \frac{15}{8}x$
- $f_3(x) = -\frac{5}{16}x^7 + \frac{21}{16}x^5 - \frac{35}{16}x^3 + \frac{35}{16}x$
- $f_4(x) = \frac{35}{128}x^9 - \frac{180}{128}x^7 + \frac{378}{128}x^5 - \frac{420}{128}x^3 + \frac{315}{128}x$

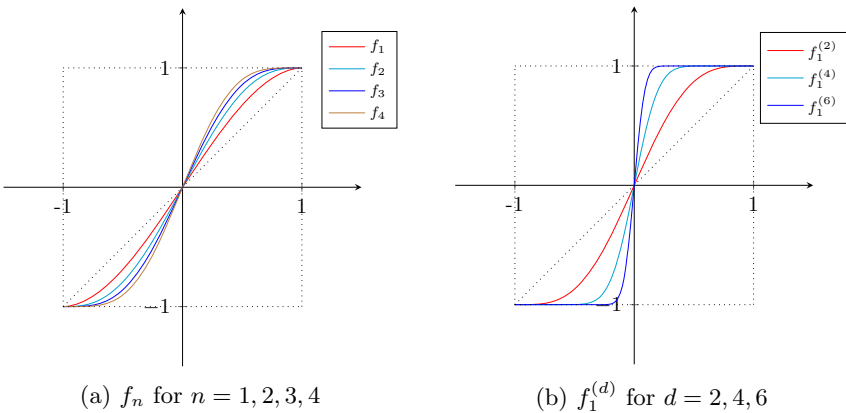


Fig. 3. Illustration of $f_n^{(d)}$

Since $\binom{2i}{i} = 2 \cdot \binom{2i-1}{i-1}$ is divisible by 2 for $i \geq 1$, every coefficient of f_n can be represented as $m/2^{2n-1}$ for $m \in \mathbb{Z}$ (Fig. 3).

Size of the Constant c_n . The constant c_n takes an important role on the convergence of $f_n^{(d)}$ (on d) to the sign function. Informally, since the coefficient of x term is exactly c_n , we can regard f_n as $f_n(x) \simeq c_n \cdot x$ for small $x > 0$, and then it holds that $1 - f_n(x) \simeq 1 - c_n \cdot x \simeq (1 - x)^{c_n}$. In the next subsection, we will present a rigorous proof of the inequality $1 - f_n(x) \leq (1 - x)^{c_n}$ for $0 < x < 1$. (see Sect. 3.2). From a simple computation, we obtain c_n as a linear summation of binomial coefficients

$$c_n = \sum_{i=0}^n \frac{1}{4^i} \binom{2i}{i},$$

which is simplified by the following lemma.

Lemma 2. *It holds that $c_n = \sum_{i=0}^n \frac{1}{4^i} \binom{2i}{i} = \frac{2n+1}{4^n} \binom{2n}{n}$.*

Proof. We prove the statement by induction. It is easy to check for $n = 1$. Assume that $c_n = \frac{2n+1}{4^n} \binom{2n}{n}$ for some $n \geq 1$. Then, it holds that

$$\begin{aligned} c_{n+1} &= c_n + \frac{1}{4^{n+1}} \binom{2n+2}{n+1} = \frac{1}{4^{n+1}} \cdot \left(\frac{2 \cdot (2n+2)!}{(n+1)!n!} + \frac{(2n+2)!}{(n+1)!(n+1)!} \right) \\ &= \frac{2n+3}{4^{n+1}} \binom{2n+2}{n+1}. \end{aligned}$$

Therefore, the lemma is proved by induction. □

To measure the size of c_n , we apply Wallis’s formula [35] which gives us very tight lower and upper bound:

$$\frac{1}{\sqrt{\pi}} \cdot \frac{2n+1}{\sqrt{n+\frac{1}{2}}} < \frac{2n+1}{4^n} \binom{2n}{n} < \frac{1}{\sqrt{\pi}} \cdot \frac{2n+1}{\sqrt{n}}.$$

From the inequality, we can check that $c_n = \Theta(\sqrt{n})$, which diverges as $n \rightarrow \infty$.

Remark 1. Our method can be naturally generalized to the composite polynomial approximation on *step functions*. For example, if we substitute Prop III by $f'(x) = cx^{2m}(1-x^2)^n$ for $m, n \geq 1$, then $f^{(d)}$ would get close to a step function F (as d increases) such that $F(x) = -1$ if $x \in [-1, -t]$, $F(x) = 0$ if $x \in [-t, t]$ and $F(x) = 1$ if $x \in (t, 1]$, for some $0 < t < 1$ as d increases.

3.2 Analysis on the Convergence of $f_n^{(d)}$

In this subsection, we analyze the convergence of $f_n^{(d)}$ to the sign function as d increases. To be precise, we give a lower bound of d which makes $f_n^{(d)}$ (α, ϵ)-close to the sign function. The following lemma gives a nice upper bound on $1 - f_n(x)$, which is even tighter than the Bernoulli’s inequality [41]: This well-known inequality implies $1 - c_n x \leq (1 - x)^{c_n}$, but since $1 - c_n x \leq 1 - f_n(x)$ we cannot directly obtain the upper bound of $1 - f_n(x)$ from this inequality.

Lemma 3. *It holds that $0 \leq 1 - f_n(x) \leq (1 - x)^{c_n}$ for $x \in [0, 1]$.*

Proof. It is trivial that $f_n(x) \leq f_n(1) = 1$ for $x \in [0, 1]$. We will prove $G(x) := (1 - x)^{c_n} - (1 - f_n(x)) \geq 0$ for $x \in [0, 1]$ by showing

1. $G(0) = G(1) = 0$,
2. there exists $x_0 \in (0, 1)$ s.t. $G(x_0) > 0$,
3. there exists a unique $y_0 \in (0, 1)$ s.t. $G'(y_0) = 0$.

We first check why these three conditions derive the result $G(x) \geq 0$. Assume that there exists $x_1 \in (0, 1)$ such that $G(x_1) < 0$. Since G is continuous, there exists a root x_2 of G between x_0 and x_1 . Then by the mean value theorem, there exist $y_1 \in (0, x_2)$ and $y_2 \in (x_2, 1)$ satisfying $G'(y_1) = G'(y_2) = 0$, which contradicts to the third condition.

Now we prove the three conditions. The first condition is trivial. To show the second condition, we observe $G(0) = 0$, $G'(0) = 0$ and $G''(0) > 0$ which can be easily checked. Since G'' is continuous, $G'(0) = 0$ and $G''(0) > 0$ imply that $G'(x) > 0$ for $x \in (0, \delta)$ for some $\delta > 0$. Combining with $G(0) = 0$, we obtain $G(x) > 0$ for $x \in (0, \delta)$ which implies the second condition.

To show the uniqueness, let $G'(x) = c_n(1 - x^2)^n - c_n(1 - x)^{c_n - 1} = 0$. Then it holds that $(1 - x)^{n - c_n + 1} \cdot (1 + x)^n = 1$ for $x \in (0, 1)$ which is equivalent to

$$\frac{\log(1 + x)}{\log(1 - x)} = -\frac{n - c_n + 1}{n}.$$

Since $\log(1 + x)/\log(1 - x)$ is a strictly increasing function, there should exist a unique $y_0 \in (0, 1)$ satisfying the equation which implies $G'(y_0) = 0$. □

We give another inequality on $1 - f_n(x)$ which is tighter than the inequality in the previous lemma when x is close to 1.

Lemma 4. *It holds that $0 \leq 1 - f_n(x) \leq 2^n \cdot (1 - x)^{n+1}$ for $x \in [0, 1]$.*

Proof. Let $y = 1 - x$, and set

$$H(y) = \frac{c_n \cdot 2^n}{n + 1} \cdot y^{n+1} - (1 - f_n(1 - y)).$$

Then $H'(y) = c_n \cdot 2^n \cdot y^n - f'_n(1 - y) = c_n \cdot 2^n \cdot y^n - c_n \cdot y^n(2 - y)^n \geq 0$ for $y \in [0, 1]$. Since $H(0) = 0$, it holds that $H(y) \geq 0$. Therefore, we obtain

$$1 - f_n(x) \leq \frac{c_n \cdot 2^n}{n + 1} \cdot (1 - x)^{n+1} \leq 2^n \cdot (1 - x)^{n+1}$$

for $x \in [0, 1]$, where the second inequality comes from $c_n < n + 1$. □

Now we obtain the theorem on the convergence of $f_n^{(d)}$ to the sign function.

Theorem 3 (Convergence of $f_n^{(d)}$). *If $d \geq \frac{1}{\log c_n} \cdot \log(1/\epsilon) + \frac{1}{\log(n+1)} \cdot \log(\alpha - 1) + O(1)$, then $f_n^{(d)}(x)$ is an (α, ϵ) -close polynomial to $\text{sgn}(x)$ over $[-1, 1]$.*

Proof. Since f_n is an odd function, it suffices to consider the case that the input x is non-negative. We analyze the lower bound of d for the convergence of $f_n^{(d)}$ by applying Lemma 3 and Lemma 4. Note that Lemma 3 is tighter than Lemma 4 if x is close to 0 while the reverse holds if x is close to 1. To this end, to obtain a tight lower bound of d , our analysis is divided into two steps:

Step 1. Since f_n is an odd function, it suffices to consider the case $x \in [\epsilon, 1]$ instead of $[-1, -\epsilon] \cup [\epsilon, 1]$. Let $d_\epsilon = \left\lceil \frac{1}{\log(c_n)} \cdot \log \left(\log \left(\frac{1}{\tau} \right) / \epsilon \right) \right\rceil$ for some constant $0 < \tau < 1$. Then applying Lemma 3, we obtain following inequality for $x \in [\epsilon, 1]$.

$$\begin{aligned} 1 - f_n^{(d_\epsilon)}(x) &\leq (1 - x)^{c_n^{d_\epsilon}} \\ &\leq (1 - \epsilon)^{\log(\frac{1}{\tau})/\epsilon} < \left(\frac{1}{e}\right)^{\log(\frac{1}{\tau})} < \tau. \end{aligned}$$

Step 2. Now let $d_\alpha = \left\lceil \frac{1}{\log(n+1)} \cdot \log \left((\alpha - 1) / \log \left(\frac{1}{2\tau} \right) \right) \right\rceil$. Applying previous result and Lemma 4, we obtain following inequality for $x \in [\epsilon, 1]$.

$$\begin{aligned} 2 \cdot \left(1 - f_n^{(d_\epsilon + d_\alpha)}(x) \right) &\leq \left(2 \cdot \left(1 - f_n^{(d_\epsilon)}(x) \right) \right)^{(n+1)^{d_\alpha}} \\ &\leq (2\tau)^{(n+1)^{d_\alpha}} \leq (2\tau)^{(\alpha-1)/\log(\frac{1}{2\tau})} = 2^{-\alpha+1}. \end{aligned}$$

Therefore, if $d \geq d_\epsilon + d_\alpha$, we obtain $1 - f_n^{(d)}(x) \leq 2^{-\alpha}$ for $x \in [\epsilon, 1]$.

Note that the choice of the constant τ is independent to ϵ and α . When $\tau = 1/4$, then we get $d_\epsilon + d_\alpha = \frac{1}{\log(c_n)} \cdot \log(1/\epsilon) + \frac{1}{\log(n+1)} \cdot \log(\alpha - 1) + \frac{1}{\log(c_n)} + O(1)$. Since $\frac{1}{\log(c_n)} \leq 2$, the theorem is finally proved. \square

Remark 2. In Appendix D, we also described the *erroneous version* of the convergence of $f_n^{(d)}$ considering the approximate error induced by HEAAN evaluation.

3.3 New Comparison Algorithm NewComp

Now we introduce our new comparison algorithm based on the previous composite function approximation (Theorem 3) of the sign function. From the identity $\text{comp}(a, b) = (\text{sgn}(a - b) + 1)/2$ and approximation $f_n^{(d)}(x) \simeq \text{sgn}(x)$, we get

$$\text{comp}(a, b) \simeq \frac{f_n^{(d)}(a - b) + 1}{2},$$

which results in our new comparison algorithm **NewComp** (Algorithm 1).

It is quite natural that the larger d gives more accurate result. Since the comparison algorithm **NewComp**($\cdot, \cdot; n, d$) is obtained from the evaluation of $f_n^{(d)}$, Theorem 3 is directly transformed into the context of **NewComp** as Corollary 1, which informs us how large d is sufficient to get the result in certain accuracy.

Algorithm 1. $\text{NewComp}(a, b; n, d)$

Input: $a, b \in [0, 1], n, d \in \mathbb{N}$

Output: An approximate value of 1 if $a > b$, 0 if $a < b$ and $1/2$ otherwise

```

1:  $x \leftarrow a - b$ 
2: for  $i \leftarrow 1$  to  $d$  do
3:    $x \leftarrow f_n(x)$  // compute  $f_n^{(d)}(a - b)$ 
4: end for
5: return  $(x + 1)/2$ 

```

Corollary 1. *If $d \geq \frac{1}{\log c_n} \cdot \log(1/\epsilon) + \frac{1}{\log(n+1)} \cdot \log(\alpha - 2) + O(1)$, then the error of the output of $\text{NewComp}(a, b; n, d)$ compared to the true value is bounded by $2^{-\alpha}$ for any $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$.*

Remark 3. One can substitute non-integer scalar multiplications in the evaluation of f_n with integer scalar multiplications by linearly transforming f_n to an integer-coefficient polynomial h_n as

$$\begin{aligned}
 h_n(x) &:= \frac{f_n(2x - 1) + 1}{2} = \sum_{i=0}^n \frac{1}{4^i} \cdot \binom{2i}{i} \cdot (2x - 1) \cdot (4x - 4x^2)^i \\
 &= \sum_{i=0}^n \binom{2i}{i} \cdot (2x - 1) \cdot (x - x^2)^i.
 \end{aligned}$$

Note that it is easily proved that $h_n^{(d)}(x) = \frac{f^{(d)}(2x-1)+1}{2}$ by induction, so we can express the comparison functions as

$$\text{comp}(a, b) \simeq \frac{f_n^{(d)}(a - b) + 1}{2} = h_n^{(d)}\left(\frac{(a - b) + 1}{2}\right).$$

Therefore, Algorithm 1 can be naturally converted into the context of h_n which does not require non-integer scalar multiplications that consume level in HE.

3.4 Computational Complexity of NewComp and Its Asymptotic Optimality

In this subsection, we analyze the computational complexity of our new comparison method, and compare the result with the previous methods. Note that the (multiplicative) computational complexity of $\text{NewComp}(\cdot, \cdot; n, d)$ equals to that of evaluating $f_n^{(d)}$, so it suffices to focus on this composite polynomial.

For each $n \geq 1$, let C_n be the required number of multiplications (hence the computational complexity) of f_n using some polynomial evaluation algorithm, and denote the lower bound of d in Theorem 3 by $d_n := \frac{1}{\log c_n} \cdot \log(1/\epsilon) + \frac{1}{\log(n+1)} \cdot \log(\alpha - 1) + O(1)$. Then the total computational complexity of $f_n^{(d_n)}$ is $TC_n := d_n \cdot C_n$ which varies on the choice of n . When n becomes larger, then

d_n becomes smaller but C_n becomes larger. Namely, there is a trade-off between d_n and C_n , so we need to find the best choice of n which minimizes the total computational complexity TC_n of $f_n^{(d_n)}$.

Table 1. Depth/Computational complexity of f_n and $f_n^{(d_n)}$

n	D_n	C_n	d_n	TD_n	TC_n
1	2	2	$L(1.71, 1)$	$L(3.42, 2)$	$L(3.42, 2)$
2	3	3	$L(1.1, 0.63)$	$L(3.3, 1.89)$	$L(3.3, 1.89)$
3	3	4	$L(0.89, 0.5)$	$L(2.67, 1.5)$	$L(3.56, 2)$
4	4	4	$L(0.77, 0.43)$	$L(3.08, 1.72)$	$L(3.08, 1.72)$
5	4	5	$L(0.7, 0.39)$	$L(2.8, 1.56)$	$L(3.5, 2.45)$
6	4	6	$L(0.64, 0.36)$	$L(2.56, 1.44)$	$L(3.84, 2.16)$
7	4	7	$L(0.61, 0.33)$	$L(2.44, 1.32)$	$L(4.27, 2.31)$

We assume that each polynomial f_n is computed by the Paterson-Stockmeyer method [43] which achieves an optimal computational complexity upto constant. Then, the depth is $D_n := \log(\deg f_n) + O(1) = \log n + O(1)$, and the computational complexity is $C_n := \Theta(\sqrt{\deg f_n}) = \Theta(\sqrt{n})^2$. The total depth of $f_n^{(d_n)}$ is $TD_n := d_n \cdot D_n = L\left(\frac{\log n + O(1)}{\log c_n}, \frac{\log n + O(1)}{\log(n+1)}\right)$ (see Sect. 2.1 for L notation). Since $c_n = \Theta(\sqrt{n})$ by Lemma 2, the total depth TD_n gets close to $L(2, 1)$ as n increases³. On the other hand, the total computational complexity of $f_n^{(d_n)}$ is

$$TC_n := d_n \cdot C_n = L\left(\frac{1}{\log c_n} \cdot \Theta(\sqrt{n}), \frac{1}{\log(n+1)} \cdot \Theta(\sqrt{n})\right),$$

which diverges as n increases, contrary to the total depth TD_n . Therefore, the optimal choice of n which minimize the total complexity TC_n exists. The exact number of multiplications C_n of f_n and the exact value of TC_n for small n 's are described in Table 1. From simple computations, we can check that $n = 4$ gives the minimal computational complexity TC_4 .

Asymptotic Optimality. As described in Sect. 2.2, the minimal degree of an (α, ϵ) -close approximate polynomial of the sign function over $[-1, 1]$ is $\Theta(\alpha/\epsilon)$. Since the sign function and the comparison function are equivalent, this implies that any comparison algorithm on inputs $a, b \in [0, 1]$ whose output is within $2^{-\alpha}$ error when $|a - b| \geq \epsilon$ requires at least $\Theta(\log \alpha) + \Theta(\log(1/\epsilon))$ complexity. As described above, the computational complexity of $\text{NewComp}(\cdot, \cdot; n, d_n)$ is

² The complexity notations in D_n and C_n only depend on n , not α and ϵ .

³ It does *not* mean the ‘‘convergence’’ to $L(2, 1)$ as $n \rightarrow \infty$, since the equation $TD_n = L\left(\frac{\log n + O(1)}{\log c_n}, \frac{\log n + O(1)}{\log(n+1)}\right)$ only holds when $n = O(1)$ with respect to α and $1/\epsilon$.

$\Theta(\log \alpha) + \Theta(\log(1/\epsilon))$ for each n . Therefore, our method achieves an *optimality in asymptotic computational complexity* upto constant, while the previous method [12] only achieves quasi-optimality with an additional $\log \alpha$ factor.

For several settings of α and ϵ , we compare the computational complexity of our method to the minimax approximation and the method in [12] as Table 2.

Table 2. Asymptotic computational complexity for each comparison method

Parameters	Minimax approx	[12] Method	Our method
$\log(1/\epsilon) = \Theta(1)$	$\Theta(\sqrt{\alpha})$	$\Theta(\log^2 \alpha)$	$\Theta(\log \alpha)$
$\log(1/\epsilon) = \Theta(\alpha)$	$\Theta(\sqrt{\alpha} \cdot 2^{\alpha/2})$	$\Theta(\alpha \cdot \log \alpha)$	$\Theta(\alpha)$
$\log(1/\epsilon) = 2^\alpha$	$\Theta(\sqrt{\alpha} \cdot 2^{2^{\alpha-1}})$	$\Theta(\alpha \cdot 2^\alpha)$	$\Theta(2^\alpha)$

3.5 Heuristic Methodology of Convergence Acceleration

In this subsection, we introduce a heuristic methodology to reduce the *constants* a and b in $L(a, b)$ of the computational complexity TC_n , which accelerates `NewComp` in practice.

The intuition of our acceleration method can be found in the proof of Theorem 3. The proof is divided into two steps: Step 1 is to make $f_n^{(d_\epsilon)}([\epsilon, 1]) \subseteq [1 - \tau, 1]$ for some constant $0 < \tau < 1$ (applying Lemma 3), and Step 2 is to make $f_n^{(d_\alpha)}([1 - \tau, 1]) \subseteq [1 - 2^{-\alpha}, 1]$ (applying Lemma 4). Our key observation is that we can accelerate Step 1 by using another function g rather than f_n . The convergence of $f_n^{(d)}$ ($1 \leq d \leq d_\epsilon$) in Step 1 mainly depends on the constant c_n , the derivative of f_n at zero. Therefore, we may expect that the required number of polynomial compositions d_ϵ in Step 1 can be reduced if we substitute f_n by some other odd polynomial g which satisfies $g'(0) > f'_n(0)$.

However, we cannot take any g with large derivative at 0, since the range of $g^{(d)}$ over the domain $[\epsilon, 1]$ must be contained in $[1 - \tau, 1]$ when d is large enough. In particular, the polynomial g must satisfy following properties (compare it with the Core Properties of f in Sect. 3.1):

- Prop I. $g(-x) = -g(x)$ (Origin Symmetry)
- Prop IV. $\exists 0 < \delta < 1$ s.t. $x < g(x) \leq 1$ for all $x \in (0, \delta]$, (Toward $[1 - \tau, 1]$)
 and $g([\delta, 1]) \subseteq [1 - \tau, 1]$ (Keep in $[1 - \tau, 1]$)

For each g , we denote the minimal δ in Prop IV by δ_0 in the rest of paper.

Note that Prop IV is necessary to make $g^{(d)}(x) \in [1 - \tau, 1]$ for $x \in [\epsilon, 1]$ when $d \geq d_0$ for some sufficiently large $d_0 > 0$. Intuitively, among all g of the same degree satisfying above properties, a smaller d is required for $g^{(d)}([\epsilon, 1]) \subseteq [1 - \tau, 1]$ if g satisfies Prop IV with smaller δ_0 and has bigger value on the interval $(0, \delta_0)$ (hence $g'(0)$ is bigger).

We introduce a novel algorithm (Algorithm 2) which outputs a degree- $(2n+1)$ polynomial denoted by $g_{n,\tau}$ having *minimal* δ_0 of Prop IV among all degree- $(2n+1)$ polynomials satisfying Prop I & IV. In a certain condition, we can additionally show that $g_{n,\tau}(x) > g(x)$ on $x \in (0, \delta)$ (hence larger derivative at zero) for any other polynomials g satisfying Prop I & IV (see Theorem 4 and Corollary 2). It implies that $g_{n,\tau}$ is the best polynomial among all same-degree polynomials achieving our goal, i.e., $g_{n,\tau}^{(d)}([\epsilon, 1]) \subseteq [1 - \tau, 1]$ with minimal d .

Algorithm 2. FindG(n, τ)

Input: $n \geq 1, 0 < \tau < 1$

Output: A degree- $(2n+1)$ polynomial $g_{n,\tau}$ satisfying Prop I & IV with minimal δ of Prop IV.

```

1:  $g_{n,\tau} \leftarrow x$  // Initialize  $g_{n,\tau}(x) = x$ 
2: repeat
3:    $\delta_0 \leftarrow$  minimal  $\delta$  s.t.  $g_{n,\tau}([\delta, 1]) \subseteq [1 - \tau, 1]$  // Initial  $\delta_0$  is  $1 - \tau$ 
4:    $g_{min} \leftarrow$  degree- $(2n+1)$  minimax approx. poly. of  $(1 - \frac{\tau}{2}) \cdot \text{sgn}(x)$  over  $[-1, -\delta_0] \cup [\delta_0, 1]$ 
5:    $g_{n,\tau} \leftarrow g_{min}$ 
6:    $S \leftarrow \|g_{n,\tau} - (1 - \frac{\tau}{2})\|_{\infty, [\delta_0, 1]}$ 
7: until  $S == \frac{\tau}{2}$ 
8: return  $g_{n,\tau}$ 

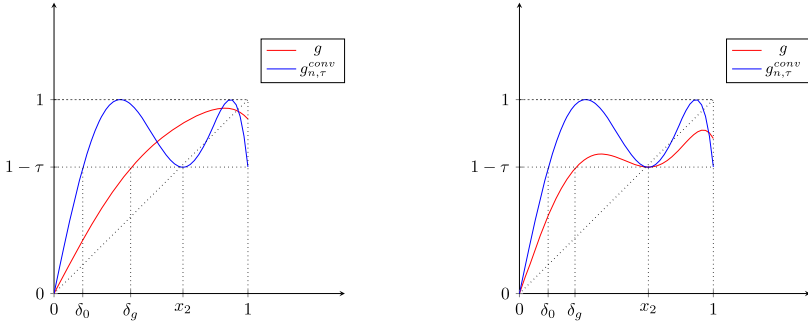
```

In Algorithm 2, the equality check $S == \frac{\tau}{2}$ on line 7 is done with a certain precision in practice (e.g., 2^{-10} or 2^{-53}). Note that S converges (increases) to $\frac{\tau}{2}$, δ_0 converges (decreases) to some $\delta_{conv} > 0$, and hence $g_{n,\tau}$ converges to some polynomial $g_{n,\tau}^{conv}$ (see Appendix B). From this, we obtain two facts: First, Algorithm 2 terminates in finite iterations given a finite precision for the equality check. Second, the algorithm output satisfies Prop I & IV⁴.

We provide a theoretical analysis on $g_{n,\tau}^{conv}$ to which $g_{n,\tau}$ converges, which we call *the ideal output polynomial* of Algorithm 2. Note that the ideal output polynomial $g_{n,\tau}^{conv}$ satisfies $\|g_{n,\tau}^{conv} - (1 - \frac{\tau}{2})\|_{\infty, [\delta_0, 1]} = \frac{\tau}{2}$. The following theorem shows the optimality of $g_{n,\tau}^{conv}$, which implies that the real output of Algorithm 2 with a certain precision is nearly optimal.

Theorem 4 (Optimality of $g_{n,\tau}^{conv}$). *The ideal output polynomial $g_{n,\tau}^{conv}$ of Algorithm 2 satisfies Prop I & IV with minimal δ_0 among all degree- $(2n+1)$ polynomials satisfying Prop I & IV. Let $x_2 > 0$ be the smallest positive x -coordinate of local minimum points of $g_{n,\tau}$ following the notation in Lemma 1 (If local minimum does not exist, set $x_2 = 1$). If $x_2 \geq 1 - \tau$, then $g_{n,\tau}(x) > g(x)$ for $x \in (0, \delta_0)$ for any other degree- $(2n+1)$ polynomial g satisfying Prop I & IV.*

⁴ In every iteration of Algorithm 2, the minimax approximate polynomial g_{min} of $(1 - \frac{\tau}{2}) \cdot \text{sgn}(x)$ over $[-1, \delta_0] \cup [\delta_0, 1]$ satisfies Prop I & IV. Prop I is trivial, and $g_{min}([\delta_0, 1]) \subset [1 - \tau, 1]$ by Lemma 1. Since $g_{min}(\delta_0) > 1 - \tau \geq \delta_0$ and g_{min} is concave & increasing in $[0, \delta_0]$, it holds that $x < g_{min}(x) < 1$ for $x \in (0, \delta_0]$.



(a) Intersections without multiplicity (b) Intersection with multiplicity at x_2

Fig. 4. Example description of intersections of g and $g_{n,\tau}^{conv}$ for $n = 3$

Proof. Let δ_{conv} be the minimal δ such that $g_{n,\tau}^{conv}([\delta, 1]) \subseteq [1 - \tau, 1]$. Assume that there exists a degree- $(2n + 1)$ polynomial g satisfying Prop I & IV with $\delta \leq \delta_{conv}$. By Prop IV, we get $\|g - (1 - \frac{\tau}{2})\|_{\infty, [\delta, 1]} \leq \frac{\tau}{2}$, and then it trivially holds that $\|g - (1 - \frac{\tau}{2})\|_{\infty, [\delta_{conv}, 1]} \leq \frac{\tau}{2} = \|g_{n,\tau}^{conv} - (1 - \frac{\tau}{2})\|_{\infty, [\delta_{conv}, 1]}$. Therefore, $g = g_{n,\tau}^{conv}$ by Lemma 1 which implies the minimality of δ_{conv} .

Now we prove the second statement. Let g be a degree- $(2n + 1)$ polynomial satisfying Prop I & IV which is distinct from $g_{n,\tau}^{conv}$, and δ_g be the minimal δ such that $g([\delta, 1]) \subseteq [1 - \tau, 1]$. From the minimality of δ_{conv} and Prop IV, it holds that $\delta_{conv} < \delta_g \leq 1 - \tau \leq x_2$. By Lemma 1, $g_{n,\tau}^{conv}$ oscillates on $[\delta_{conv}, 1]$ with 1 and $1 - \tau$ as maximum and minimum, respectively, and it has n critical points in $(\delta_{conv}, 1)$. Since $g([\delta_g, 1]) \subseteq [1 - \tau, 1]$ and $\delta_g \leq x_2$, the polynomial g intersects with $g_{n,\tau}^{conv}$ on at least n points in $[\delta_g, 1]$: when $g(x) = g_{n,\tau}^{conv}(x)$ and $g'(x) = g_{n,\tau}^{conv}'(x)$, then x is counted as two points (see Fig. 4). Now our second argument is proved as following: If $g(x) \geq g_{n,\tau}^{conv}(x)$ ⁵ on some $x \in (0, \delta_{conv}) \subset (0, \delta_g)$, then g and $g_{n,\tau}^{conv}$ intersect on at least one point in $(0, \delta_g)$ by intermediate value theorem since there exists $y \in (\delta_{conv}, \delta_g)$ such that $g(y) < 1 - \tau \leq g_{n,\tau}^{conv}(y)$ by the definition of δ_g . This leads to a contradiction since g and $g_{n,\tau}^{conv}$ intersect on $2(n + 1) + 1 = 2n + 3$ points (the factor 2 comes from the fact that both are odd polynomials) including the origin while the degree of both g and $g_{n,\tau}^{conv}$ is $2n + 1 < 2n + 3$. Therefore, $g_{n,\tau}^{conv}(x) > g(x)$ for all $x \in (0, \delta_{conv})$. \square

Corollary 2. Let $g_{n,\tau}^{conv}$ be the ideal output polynomial of Algorithm 2, and δ_0 be the corresponding minimal δ satisfying Prop IV. If $n = 1$, $(n, \tau) = (2, 0.25)$, or $(n, \tau) = (3, 0.35)$, then $\delta_0 < \delta_g$ and $g_{n,\tau}^{conv}(x) > g(x)$ on $x \in (0, \delta_0)$ for any other degree- $(2n + 1)$ polynomial g satisfying Prop I & IV.

Though $g_{n,\tau}$ is hard to be expressed in closed form contrary to f_n , we can find it with a certain precision (e.g., 2^{-10}) by running Algorithm 2 in MATLAB. For

⁵ If $g(x) = g_{n,\tau}^{conv}(x)$ on some $x \in (0, \delta_0)$, it is the point of intersection in $(0, \delta_g)$, and proof continues.

example, we provide explicit descriptions of the polynomials $g_{n,\tau}$ for $n = 1, 2, 3, 4$ and $\tau = \frac{1}{4}$ (Fig. 5). In this case, the equality check in Algorithm 2 was done with 10^{-4} precision. We omit the subscript τ of $g_{n,\tau}$ for $\tau = \frac{1}{4}$ for convenience.

- $g_1(x) = -\frac{1359}{2^{10}} \cdot x^3 + \frac{2126}{2^{10}} \cdot x$
- $g_2(x) = \frac{3796}{2^{10}} \cdot x^5 - \frac{6108}{2^{10}} \cdot x^3 + \frac{3334}{2^{10}} \cdot x$
- $g_3(x) = -\frac{12860}{2^{10}} \cdot x^7 + \frac{25614}{2^{10}} \cdot x^5 - \frac{16577}{2^{10}} \cdot x^3 + \frac{4589}{2^{10}} \cdot x$
- $g_4(x) = \frac{46623}{2^{10}} \cdot x^9 - \frac{113492}{2^{10}} \cdot x^7 + \frac{97015}{2^{10}} \cdot x^5 - \frac{34974}{2^{10}} \cdot x^3 + \frac{5850}{2^{10}} \cdot x$

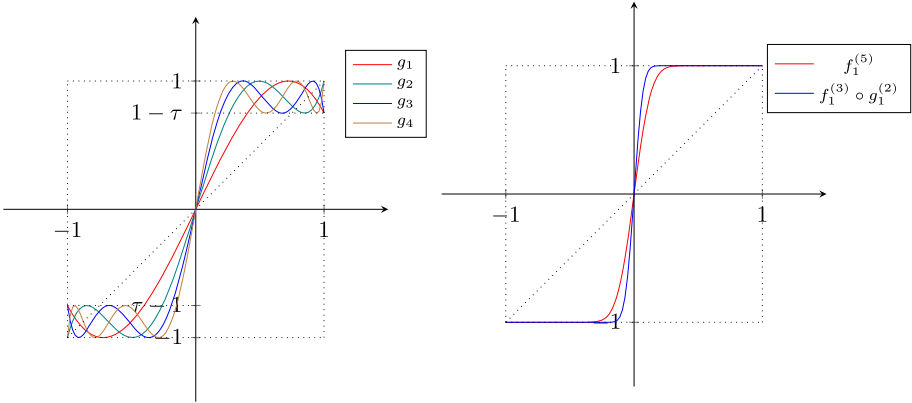


Fig. 5. Illustration of g_n and the comparison of $f_1^{(d_f+d_g)}$ and $f_1^{(d_f)} \circ g_1^{(d_g)}$

We can empirically check that g_n also satisfies the following two heuristic properties. The first property shows how large $g'_n(0)$ is when it compared to $f'_n(0)$, and the second property shows how fast $g_n(x)$ gets close to ± 1 , i.e., the g_n -version of Lemma 3.

Heuristic Properties of g_n :

1. $g'_n(0) \simeq 0.98 \cdot f'_n(0)^2$ (Hence, $\log g'_n(0) \simeq 2 \cdot \log c_n$)
2. $1 - g_n(x) \leq (1 - x)^{g'_n(0)}$ for $x \in [0, \delta_0]$ where δ_0 is the minimal δ in Prop IV

Experimental results supporting above heuristic properties are described in Appendix C. Applying these g_n polynomials, we can provide a new comparison algorithm (Algorithm 3), which is a modified version of Algorithm 1 and offers the same functionality with the reduced computational complexity and depth. We can also estimate the number of compositions d_f and d_g required for this modified algorithm to achieve a certain accuracy as Corollary 3.

Corollary 3 (With Heuristic Properties). *If $d_g \geq \frac{1}{\log g'_n(0)} \cdot \log(1/\epsilon) + O(1) = \frac{1/2+o(1)}{\log c_n} \cdot \log(1/\epsilon) + O(1)$ and $d_f \geq \frac{1}{\log n} \cdot \log(\alpha - 2) + O(1)$, then the error of the output of $\text{NewCompG}(a, b; n, d_f, d_g)$ compared to the true value is bounded by $2^{-\alpha}$ for any $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$.*

Algorithm 3. NewCompG($a, b; n, d_f, d_g$)

Input: $a, b \in [0, 1], n, d_f, d_g \in \mathbb{N}$

Output: An approximate value of 1 if $a > b$, 0 if $a < b$ and $1/2$ otherwise

```

1:  $x \leftarrow a - b$ 
2: for  $i \leftarrow 1$  to  $d_g$  do
3:    $x \leftarrow g_n(x)$  // compute  $g_n^{(d_g)}(a - b)$ 
4: end for
5: for  $i \leftarrow 1$  to  $d_f$  do
6:    $x \leftarrow f_n(x)$  // compute  $f_n^{(d_f)} \circ g_n^{(d_g)}(a - b)$ 
7: end for
8: return  $(x + 1)/2$ 

```

Proof. Following the proof of Theorem 3, it suffices to show that $1 - g_n^{(d_g)}(x) \leq \tau$ for $x \in [\epsilon, 1]$ where $\tau = 1/4$. Let $e_n := g'_n(0)$. By the second heuristic property of g_n , we obtain two inequalities: $1 - g_n^{(d)}(x) \leq (1 - x)e_n^d$ for d satisfying $g_n^{(d-1)}(x) \leq \delta_0$, and $1 - g_n^{(d)}(x) \leq \tau$ for $g_n^{(d-1)}(x) > \delta_0$. Therefore, it holds that

$$1 - g_n^{(d)}(x) \leq \max\left((1 - x)e_n^d, \tau\right)$$

for any $d > 0$. Applying $d = d_g := \left\lceil \frac{1}{\log e_n} \cdot \log\left(\log\left(\frac{1}{\tau}\right) / \epsilon\right) \right\rceil$, we finally obtain $1 - g_n^{(d_g)}(x) \leq \tau$ since $(1 - x)e_n^{d_g} \leq (1 - \epsilon)^{\log(\frac{1}{\tau})/\epsilon} < \tau$. □

The important point is that d_g is reduced as approximately *half* (applying the first heuristic property of g_n) compared to the previous case that only uses f_n to approximate the sign function. Since g_n and f_n requires same number of non-scalar multiplications, we can conclude that the computational complexity of $f_n^{(d_f)} \circ g_n^{(d_g)}$ is $L\left(\frac{a_n}{2}, b_n\right)$ where a_n and b_n are defined from $TC_n = L(a_n, b_n)$.

The total depth of $f_n^{(d_f)} \circ g_n^{(d_g)}$ is $L\left(\frac{\log n + O(1)}{2 \cdot \log c_n}, \frac{\log n + O(1)}{\log(n+1)}\right)$ which gets close to $L(1, 1)$ as n increases⁶. Note that $L(1, 1)$ is theoretically the minimal depth obtained by minimax polynomial approximation (see Sect. 2.2).

4 Application to Min/Max

As described in [12], min/max functions correspond to the absolute function as

$$\min(a, b) = \frac{a + b}{2} - \frac{|a - b|}{2} \quad \text{and} \quad \max(a, b) = \frac{a + b}{2} + \frac{|a - b|}{2}.$$

Therefore, an approximate polynomial of $|x|$ directly gives us the approximate polynomial of min/max functions. Since $|x| = x \cdot \text{sgn}(x)$, we can consider the

⁶ It does *not* mean the “convergence” to $L(1, 1)$ as $n \rightarrow \infty$, since n should be $O(1)$ with respect to α and $1/\epsilon$.

convergence of $x \cdot f_n^{(d)}(x)$ to $|x|$ as an analogue. As $\min(a, b)$ is directly computed from $\max(a, b)$, we only describe an algorithm of \max for convenience.

Contrary to $\text{sgn}(x)$, the absolute function $|x|$ is continuous so that the parameter ϵ is unnecessary. The following theorem provides the convergence rate of $x \cdot f_n^{(d)}(x)$ to $|x|$.

Theorem 5 (Convergence of $x \cdot f_n^{(d)}$). *If $d \geq \frac{1}{\log c_n} \cdot (\alpha - 1)$, then the error of $x \cdot f_n^{(d)}(x)$ compared to $|x|$ is bounded by $2^{-\alpha}$ for any $x \in [-1, 1]$.*

Proof. Since $|x| = x \cdot \text{sgn}(x)$, the error is upper bounded as

$$\left| x \cdot f_n^{(d)}(x) - |x| \right| = |x| \cdot \left| f_n^{(d)}(x) - \text{sgn}(x) \right| \leq |x| \cdot |1 - |x||^{c_n^d}.$$

Let $y = |x| \in [0, 1]$ and $k = c_n^d$, then the error upper bound is expressed as $E(y) = y \cdot (1 - y)^k$. By a simple computation, one can check that $E(y)$ has the maximal value at $y = 1/(k + 1)$. Therefore, k should satisfy

$$E\left(\frac{1}{k + 1}\right) = \frac{k^k}{(k + 1)^{k+1}} \leq 2^{-\alpha}.$$

Since $2 \leq (1 + 1/k)^k \leq e$ for $k \geq 1$, setting $k \geq 2^{\alpha-1}$ implies $d \geq \frac{1}{\log c_n} \cdot (\alpha - 1)$. □

We denote an algorithm which evaluates $\frac{a+b}{2} + \frac{a-b}{2} \cdot f_n^{(d)}(a - b)$ by **NewMax** (see Algorithm 4), and Theorem 5 is naturally transformed into the context of \min/\max as Corollary 4.

Algorithm 4. **NewMax**($a, b; n, d$)

Input: $a, b \in [0, 1], n, d \in \mathbb{N}$

Output: An approximate value of $\max(a, b)$

1: $x \leftarrow a - b, y \leftarrow \frac{a+b}{2}$

2: **for** $i \leftarrow 1$ **to** d **do**

3: $x \leftarrow f_n(x)$

// compute $f_n^{(d)}(a - b)$

4: **end for**

5: $y \leftarrow y + \frac{a-b}{2} \cdot x$

6: **return** y

Corollary 4. *If $d \geq \frac{1}{\log c_n} \cdot (\alpha - 2)$, then the error of the output of **NewMax**($a, b; n, d$) compared to the true value is bounded by $2^{-\alpha}$ for any $a, b \in [0, 1]$.*

Our Max v.s. Previous Max. In [12], Cheon et al. introduced a \max algorithm exploiting the same identity $\max(a, b) = \frac{a+b}{2} + \frac{|a-b|}{2}$, but they interpret the

absolute function as $|x| = \sqrt{x^2}$ which is different with our our interpretation $|x| = x \cdot \text{sgn}(x)$. To compute $\sqrt{(a-b)^2}$, they exploit Wilkes’s algorithm [48] denoted by $\text{Sqrt}(y; d)$ which approximately computes \sqrt{y} for $y \in [0, 1]$: Let $a_0 = y$ and $b_0 = y - 1$, and iteratively compute $a_{n+1} = a_n \left(1 - \frac{b_n}{2}\right)$ and $b_{n+1} = b_n^2 \left(\frac{b_n-3}{4}\right)$ for $0 \leq n \leq d - 1$, where the final output is a_d .

We note that the output of $\text{Sqrt}(x^2; d)$ equals to $x \cdot f_1^{(d)}(x)$, which means our max algorithm $\text{NewMax}(a, b; 1, d)$ (in the case of $n = 1$) gives the same output to the max algorithm in [12]. However, there are several significant advantages to use our max algorithm instead of the max algorithm in [12].

- $\text{Sqrt}(x^2; d)$ requires 3 multiplications including 1 square multiplication for each iteration, while $f_1(x)$ can be computed by only 2 multiplications. Therefore, $\text{NewMax}(\cdot, \cdot; 1, d_1)$ is faster than the max algorithm in [12].
- We can further optimize our max algorithm by substituting $f_1(x)$ with $f_n(x)$ for some $n > 1$. As an analogue of Sect. 3.4, we can select an optimal n which minimizes $d \cdot C_n$ where $d = \frac{1}{\log c_n} \cdot (\alpha - 2)$, where $n = 4$ is optimal.
- Applying the approximate HE scheme HEAAN [10, 11], the max algorithm in [12] is unstable when two inputs a and b are too close. To be precise, if the input $(a - b)^2$ is close to zero and even smaller than an error accompanied by HEAAN, then the input attached with the error can be a negative value. However, the output of $\text{Sqrt}(y; d)$ for $y < 0$ diverges as d increases. In contrary, $f_n^{(d)}$ is stable over the interval $[-1, 1]$, so our max algorithm still works well even if two inputs are very close.

Applying $\{g_n\}_{n \geq 1}$ to Max. As a construction of NewCompG , we can also apply the family of polynomials $\{g_n\}_{n \geq 1}$ with heuristic properties to accelerate our NewMax algorithm. We denote an algorithm which evaluates $\frac{a+b}{2} + \frac{a-b}{2} \cdot f^{(d_f)} \circ g^{(d_g)}(a-b)$ by $\text{NewMaxG}(a, b; n, d_f, d_g)$. Applying $\epsilon = 2^{-\alpha}$ to Corollary 3, one can easily obtain the following result on NewMaxG .

Corollary 5. *If $d_g \geq \frac{1}{\log g'_n(0)} \cdot \alpha + O(1)$ and $d_f \geq \frac{1}{\log n} \cdot \log(\alpha - 2) + O(1)$, then the error of the output of $\text{NewMaxG}(a, b; n, d_f, d_g)$ compared to the true value is bounded by $2^{-\alpha}$.*

5 Experimental Results

We measured the performance of our algorithms with comparison to Comp or Max of [12]. The experiments are divided into two categories: 1. Running algorithms on plain inputs, 2. Running algorithms on encrypted inputs. All experiments were conducted on Linux with Intel Xeon CPU at 2.10GHz processor with 8 threads. For experiments in an encrypted state, we used HEAAN library [11, 45].

5.1 Approximate HE Scheme HEAAN

Cheon et al. [11] proposed an HE scheme HEAAN which supports approximate computations of real/complex numbers. Let N be a power-of-two integer and

L be the bit-length of initial ciphertext modulus, and define $q_\ell = 2^\ell$ for $1 \leq \ell \leq L$. For $R = \mathbb{Z}[X]/(X^N + 1)$ and $R_q := R/qR$, let χ_{key} , χ_{err} and χ_{enc} be distributions over R . A (field) isomorphism $\tau : \mathbb{R}[X]/(X^N + 1) \rightarrow \mathbb{C}^{N/2}$ is applied for encoding/decoding of plaintexts.

- **KeyGen**(N, L, D).
 - Sample $s \leftarrow \chi_{\text{key}}$. Set the secret key as $\text{sk} \leftarrow (1, s)$.
 - Sample $a \leftarrow U(R_{q_L})$ and $e \leftarrow \chi_{\text{err}}$. Set $\text{pk} \leftarrow (-a \cdot s + e, a) \in R_{q_L}^2$.
 - Sample $a' \leftarrow U(R_{q_L^2})$ and $e' \leftarrow \chi_{\text{err}}$, and set $\text{evk} \leftarrow (b' = -a' \cdot s + e' + q_L \cdot s^2, a') \in R_{q_L^2}^2$.
- **Enc**_{pk}($\mathbf{m}; \Delta$).
 - For a plaintext $\mathbf{m} = (m_0, \dots, m_{N/2-1})$ in $\mathbb{C}^{N/2}$ and a scaling factor $\Delta = 2^p > 0$, compute a polynomial $\mathbf{m} \leftarrow \lfloor \Delta \cdot \tau^{-1}(\mathbf{m}) \rfloor \in R$
 - Sample $v \leftarrow \chi_{\text{enc}}$ and $e_0, e_1 \leftarrow \chi_{\text{err}}$. Output $\text{ct} = [v \cdot \text{pk} + (\mathbf{m} + e_0, e_1)]_{q_L}$.
- **Dec**_{sk}($\text{ct}; \Delta$).
 - For a ciphertext $\text{ct} = (c_0, c_1) \in R_{q_\ell}^2$, compute $\mathbf{m}' = [c_0 + c_1 \cdot s]_{q_\ell}$.
 - Output a plaintext vector $\mathbf{m}' = \Delta^{-1} \cdot \tau(\mathbf{m}') \in \mathbb{C}^{N/2}$.
- **Add**(ct, ct'). For $\text{ct}, \text{ct}' \in R_{q_\ell}^2$, output $\text{ct}_{\text{add}} \leftarrow [\text{ct} + \text{ct}']_{q_\ell}$.
- **Mult**_{evk}(ct, ct'). For $\text{ct} = (c_0, c_1), \text{ct}' = (c'_0, c'_1) \in R_{q_\ell}^2$, let $(d_0, d_1, d_2) = (c_0 c'_0, c_0 c'_1 + c_1 c'_0, c_1 c'_1)$. Compute $\text{ct}'_{\text{mult}} \leftarrow [(d_0, d_1) + \lfloor q_L^{-1} \cdot d_2 \cdot \text{evk} \rfloor]_{q_\ell}$, and output $\text{ct}_{\text{mult}} \leftarrow \lfloor \lfloor \Delta^{-1} \cdot \text{ct}'_{\text{mult}} \rfloor \rfloor_{q_{\ell-p}}$.

The secret key distribution χ_{key} is set to be $\mathcal{HWT}_N(256)$, which uniformly samples an element with ternary coefficients in R that has 256 non-zero coefficients.

5.2 Parameter Selection

We have two parameters α and ϵ which measure the quality of our comparison algorithms. In our experiments, we set $\epsilon = 2^{-\alpha}$, which is the case expecting that input and output of algorithms have the same precision bits.

HEAAN Parameters. We fix the dimension $N = 2^{17}$, then we can set the initial ciphertext modulus q_L upto 2^{1700} to achieve 128-bit security estimated by Albrecht's LWE estimator [1] (Refer to Appendix E for the script). In each experiment, we set the initial modulus such that the modulus bit after each algorithm is $\log \Delta + 10$. For example, on our comparison algorithm **NewComp**($\cdot, \cdot; n, d$), we set the initial modulus bit as

$$\log q_L = (\log \Delta \cdot \lceil \log(2n + 1) \rceil + 2n - 1) \cdot d + \log \Delta + 10.$$

Note that each coefficient of f_n is of the form $m/2^{2n-1}$ for $m \in \mathbb{Z}$ (Sect. 3.1). We progress the scalar multiplication of $m/2^{2n-1}$ in an encrypted state by multiplying m and scaling $(2n - 1)$ bits down which results in the factor $(2n - 1)$ in the above equation. In the case of **NewCompG**($\cdot, \cdot; n, d_f, d_g$), we similarly set

$$\log q_L = \log \Delta \cdot \lceil \log(2n + 1) \rceil \cdot (d_f + d_g) + (2n - 1) \cdot d_f + 10 \cdot d_g + \log \Delta + 10.$$

The bit-length of the scaling factor Δ is set to be around 40 as in [12].

Note that one can evaluate $N/2$ comparison functions simultaneously in a single homomorphic comparison. In this sense, an amortized running time of our algorithm is obtained by dividing the total running time by $N/2 = 2^{16}$.

Choice of n in $\{f_n\}_{n \geq 1}$ and $\{g_n\}_{n \geq 1}$. One should consider a different cost model other than TC_n in the case of experiments in an encrypted state. When running our algorithms with HEAAN, not only the complexity TC_n but also the depth TD_n is an important factor affecting the running time, since the computational cost of a homomorphic multiplication is different for each level. Instead of TC_n , we take another cost model $TD_n \cdot TC_n$ considering that a multiplication in R_q takes (quasi-)linear time with respect to $\log q$. Under the setting $\epsilon = 2^{-\alpha}$, one can check by simple computation that $n = 4$ also minimizes $TD_n \cdot TC_n$ as well as TC_n , and we used f_n and g_n with $n = 4$ for the experiments.

5.3 Performance of NewComp and NewCompG

We compare the performance of our new comparison algorithms **NewComp** and **NewCompG** with the previous comparison algorithm **Comp** proposed in [12]. The following experimental results show that **NewComp** is much faster than **Comp** in practice, and applying g_n polynomials (**NewCompG**) substantially improves the performance of **NewComp**.

Plain State Experiment. For “plain inputs” $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon = 2^{-\alpha}$, we measured the required computational complexity and depth of each comparison algorithm to obtain an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error. The parameters d, d_f and d_g are chosen as the lower bounds described in Corollary 1 and Corollary 3, and we checked that these theoretical lower bounds are indeed very close to those obtained experimentally.

From Fig. 6, we can see that **NewComp** requires much less depth and complexity than **Comp**, and those of **NewCompG** are even smaller. Note that the gap between these algorithms in terms of both depth and complexity grows up as α increases. For example, when $\alpha = 8$, the required complexity is $\times 3$ – 4 less in **NewComp** and **NewCompG**; when $\alpha = 32$, it is over $\times 7$ less in **NewCompG**.

Table 3. Running time (amortized running time) of **Comp**, **NewComp** and **NewCompG** on HEAAN for various α and $\epsilon = 2^{-\alpha}$; an asterisk (*) means that the parameter for HEAAN does not achieve 128-bit security due to large $\log q_L \geq 1700$.

α	Comp	NewComp	NewCompG
8	238 s (3.63 ms)*	59 s (0.90 ms)	31 s (0.47 ms)
12	572 s (8.73 ms)*	93 s (1.42 ms)	47 s (0.72 ms)
16	1429 s (21.8 ms)*	151 s (2.30 ms)*	80 s (1.22 ms)
20	2790 s (42.6 ms)*	285 s (4.35 ms)*	94 s (1.43 ms)*

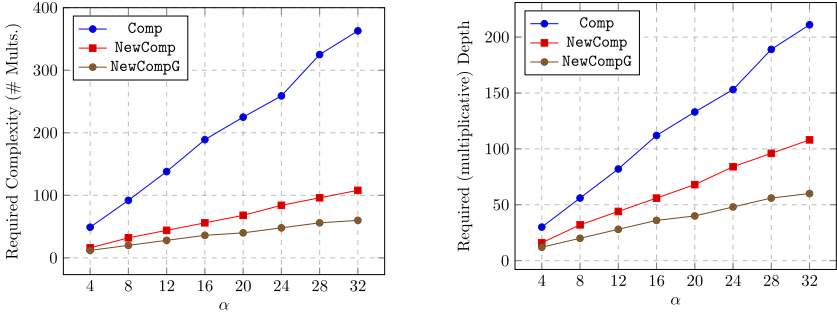


Fig. 6. Comp, NewComp and NewCompG on various α with $\epsilon = 2^{-\alpha}$ in a plain state

Encrypted State Experiment. We also measured the performance of our algorithms which output an approximate value of $\text{comp}(a, b)$ within $2^{-\alpha}$ error for “encrypted inputs” $a, b \in [0, 1]$ satisfying $|a - b| \geq \epsilon$. Note that parameters d , d_f and d_g are chosen as the lower bounds in Corollary 1 and 3. We checked through 100 experiments that our algorithms with chosen parameters give accurate results in spite of errors accompanied by HEAAN.

In Table 3, we can see the running time (and amortized running time) of our algorithms NewComp, NewCompG, and that of Comp [12] for various α . Note that our new algorithms NewComp and NewCompG provide outstanding performance in terms of amortized running time: NewComp takes 0.9 ms for 8-bit comparison, and NewCompG only takes about 1 ms to compare up to 20-bit inputs. It is a significant improvement over the previous algorithm Comp. For example, NewCompG is about $\times 8$ faster than Comp when $\alpha = 8$, about $\times 18$ faster when $\alpha = 16$, and the ratio increases as α increases.

Note that the required depth of Comp is much larger than that of our algorithms as described in Fig. 6. Consequently, to run Comp for $\alpha \geq 10$ in an encrypted state with 128-bit security, one must increase the HEAAN parameter from $N = 2^{17}$ to $N = 2^{18}$, or use bootstrapping techniques [10], both of which yields more than twice performance degradation, especially in total running time.

5.4 Performance of NewMax and NewMaxG

We also compared the performance of NewMax and NewMaxG in an encrypted state to that of the max algorithm Max in the previous work [12]. The parameters d , d_f and d_g were chosen from the theoretical lower bounds described in Corollary 4 and Corollary 5, and were confirmed that they are very close to those obtained experimentally. In Fig. 7, we can see the running time of our new algorithms NewMax, NewMaxG, and that of Max in [12]. Our algorithms improve the Max considerably in running time (and depth), and the gap increases for larger α : when $\alpha = 8$, our NewMax and NewMaxG algorithms are $\times 1.6$ and $\times 2$ faster than Max, respectively; when $\alpha = 20$, our NewMaxG algorithm is $\times 4.5$ faster than Max.

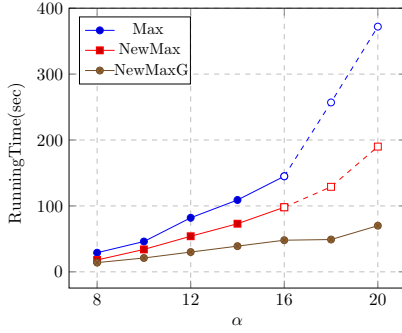


Fig. 7. Running Time of Max, NewMax and NewMaxG on HEAAN for various α . Hollow marker implies that the parameter does not achieve 128-bit security due to $\log q_L \geq 1700$.

Acknowledgement. We thank Kyoohyung Han for useful discussions in the early stage of this work, anonymous reviewers of Eurocrypt 2020 for suggesting us to investigate the line of work in numerical analysis, and those of Asiacrypt 2020 for valuable comments. This work was supported by the Institute for Information & Communications Technology Promotion (IITP) Grant through the Korean Government (MSIT), (Development and Library Implementation of Fully Homomorphic Machine Learning Algorithms supporting Neural Network Learning over Encrypted Data), under Grant 2020-0-00840.

A Derivation of f_n from Core Properties

Given $f_n(x) = c_n \int_0^x (1 - s^2)^n ds$, we use substitution $s = \sin t$ to get $\frac{f_n(x)}{c_n} = \int_0^{\sin^{-1} x} \cos^{2n+1} t dt$. Applying the following identity (which holds for any $m \geq 2$)

$$\int_0^x \cos^m t dt = \frac{1}{m} \cdot \cos^{m-1} x \cdot \sin x + \frac{m-1}{m} \cdot \int_0^x \cos^{m-2} t dt,$$

we obtain

$$\frac{f_n(x)}{c_n} = \frac{1}{2n+1} (1-x^2)^n x + \frac{2n}{2n+1} \frac{f_{n-1}(x)}{c_{n-1}}$$

for $n \geq 2$, and $\frac{f_1(x)}{c_1} = \frac{1}{3}(1-x^2)x + \frac{2}{3} \cdot x$. By induction, we can obtain f_n as

$$\frac{f_n(x)}{c_n} = \frac{1}{2n+2} \cdot \sum_{i=0}^n \prod_{k=i}^n \frac{2k+2}{2k+1} \cdot (1-x^2)^i x \tag{1}$$

Now, since $f_n(1) = 1$, evaluating above equation at $x = 1$ gives,

$$c_n = \prod_{k=1}^n \frac{2k+1}{2k} = \frac{1}{4^n} \binom{2n}{n} (2n+1).$$

Substituting this c_n into Eq. (1) and arranging, we get

$$f_n(x) = \sum_{i=0}^n \frac{1}{4^i} \cdot \binom{2i}{i} \cdot x(1-x^2)^i.$$

B Convergence of δ_0 , S and $g_{n,\tau}$

It is trivial that $S \leq \frac{\tau}{2}$. Let us denote S , δ_0 and $g_{n,\tau}$ updated in the i -th iteration by S_i , $\delta_{0,i}$ and $g_{n,\tau,i}$ respectively. Assume that $S_i < \frac{\tau}{2}$ for some $i \geq 1$. Then it holds that $g_{n,\tau,i}(x) \geq (1 - \frac{\tau}{2}) - S_i > 1 - \tau$ for $x \in [\delta_{0,i}, 1]$. Therefore, $\delta_{0,i+1}$ should be smaller than $\delta_{0,i}$, and hence S_{i+1} is larger than S_i . Since $\delta_{0,i}$ has a lower bound 0, $\delta_{0,i}$ converges to some constant $\delta_{conv} > 0$ as i increases. Hence, $g_{n,\tau,i}$ converges to some $g_{n,\tau}^{conv}$, and S_i converges to some $S_{conv} \leq \frac{\tau}{2}$.

Now, assume that $S_{conv} < \frac{\tau}{2}$ and let $\rho = \frac{\tau}{2} - S_{conv} > 0$. Since $\delta_{0,i}$ converges (and decreases) to δ_{conv} , there exists some $i \geq 1$ such that $\delta_{0,i} < \frac{1-\tau+\rho}{1-\tau} \cdot \delta_{conv}$. Note that $g_{n,\tau,i}$ is concave in $[0, \delta_{0,i}]$ as noted in Sect. 2.2. Therefore, it holds that $\frac{g_{n,\tau,i}(\delta_{0,i}) - (1-\tau)}{\delta_{0,i} - \delta_{0,i+1}} < \frac{g_{n,\tau,i}(\delta_{0,i})}{\delta_{0,i}}$ where $g_{n,\tau,i}(\delta_{0,i+1}) = 1 - \tau$. Since $g_{n,\tau,i}(\delta_{0,i}) - (1 - \tau) \geq \rho$, we obtain

$$\begin{aligned} \delta_{0,i} - \delta_{0,i+1} &> \frac{g_{n,\tau,i}(\delta_{0,i}) - (1 - \tau)}{g_{n,\tau,i}(\delta_{0,i})} \delta_{0,i} = \delta_{0,i} - \frac{1 - \tau}{g_{n,\tau,i}(\delta_{0,i})} \delta_{0,i} \\ &\geq \delta_{0,i} - \frac{1 - \tau}{1 - \tau + \rho} \delta_{0,i} = \frac{\rho}{1 - \tau + \rho} \delta_{0,i}. \end{aligned}$$

Hence, we get $\delta_{0,i} > \frac{1-\tau+\rho}{1-\tau} \cdot \delta_{0,i+1} \geq \frac{1-\tau+\rho}{1-\tau} \cdot \delta_{conv}$, which is a contradiction.

C Heuristic Properties on g_n

We provide experimental results validating the heuristic properties in Sect. 3.5:

1. $g'_n(0) \simeq 0.98 \cdot f'_n(0)^2$ (Hence, $\log g'_n(0) \simeq 2 \cdot \log c_n$)
2. $1 - g_n(x) \leq (1 - x)^{g'_n(0)}$ for $x \in [0, \delta_0]$ where δ_0 is the minimal δ in Prop IV

On the First Heuristic. Using MATLAB, we computed $g'_n(0)$ and compared it with $f_n'^2(0)$ derived from Lemma 2. See Fig. 8 for $1 \leq n \leq 20$.

On the Second Heuristic. Let $G_n(x) := 1 - (1 - x)^{g'_n(0)}$, then we can experimentally check that $G_n(x) \leq g_n(x)$ when $x \in (0, \delta_0]$, which is equivalent to $1 - g_n(x) \leq (1 - x)^{g'_n(0)}$. Let δ_1 be the largest δ such that $G_n(x) \leq g_n(x)$ for all $x \in [0, \delta]$ (see Fig. 9a). The experiment results show that $1/\delta_0 > 1/\delta_1$ which is equivalent to $\delta_0 < \delta_1$ (see Fig. 9b for $1 \leq n \leq 20$).

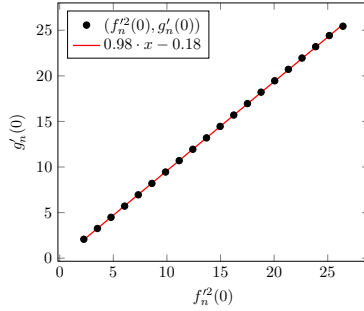


Fig. 8. $f_n^2(0)$ and $g'_n(0)$ ($R^2 = 0.9999$); $n = 1, 2, \dots, 20$ from the left to the right

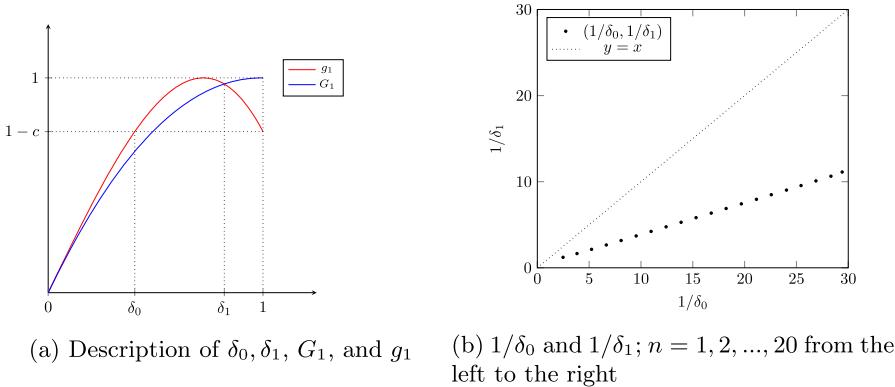


Fig. 9. Experimental evidence on $1 - g_n(x) \leq (1 - x)^{g'_n(0)}$ when $x \in (0, \delta_0]$

D Convergence of $f_n^{(d)}$ in Erroneous Case

Due to the approximate nature of HEAAN, the evaluation of f_n on an input x in an encrypted stated output an approximate value of $f_n(x)$ rather than the exact value. In this section, we analyze the convergence of $f_n^{(d)}$ considering errors induced by HEAAN evaluation, and show that the convergence is still valid in some sense under some conditions on parameters.

Let us denote by $\hat{f}_n(x)$ an approximate value of $f_n(x)$ obtained from HEAAN evaluation, i.e., $|\hat{f}_n(x) - f_n(x)| \approx 0$. For a fixed $n \geq 1$, let us assume that an approximate error $|\hat{f}_n(x) - f_n(x)|$ is bounded by $0 < B \ll 1$ (e.g., $B \approx 2^{-20}$). Then it holds that

$$|1 - \hat{f}_n(x)| \leq |1 - f_n(x)| + B.$$

Note that B can be easily controlled by changing the scaling factor Δ of HEAAN.

Now we provide some variants of Lemma 3 and Lemma 4 considering the approximation errors. To simplify the proofs, we assume that $n \geq 3$ so that $c_n > 2$.

Lemma 5. Let $B \leq \left(\frac{1}{2^{n+1}}\right)^{\frac{c_n-1}{n}} - \left(\frac{1}{2^{n+1}}\right)^{\frac{c_n}{n}}$. For $\left(\frac{c_n}{c_n-1}\right)^{c_n-1} \cdot B \leq x \leq 1 - \left(\frac{1}{2^{n+1}}\right)^{\frac{1}{n}}$, it holds that $-B < 1 - \hat{f}_n(x) < (1-x)^{c_n-1}$.

Proof. The first inequality is trivial since $\hat{f}_n(x) \leq f_n(x) + B \leq 1 + B$. For $K(x) = (1-x)^{c_n-1} - (1-x)^{c_n}$, it is easy to check that K has a unique local maximal point $\left(x_0 = \frac{1}{c_n}, K(x_0) = \frac{1}{c_n} \cdot \left(1 - \frac{1}{c_n}\right)^{c_n-1}\right)$ over $[0, 1]$ and is convex in $[0, x_0]$. As a result, for $\frac{x_0}{K(x_0)} \cdot B = \left(\frac{c_n}{c_n-1}\right)^{c_n-1} \cdot B \leq x \leq \frac{1}{c_n}$, it holds that $B \leq K(x)$. Since $B \leq \left(\frac{1}{2^{n+1}}\right)^{\frac{c_n-1}{n}} - \left(\frac{1}{2^{n+1}}\right)^{\frac{c_n}{n}} = K\left(1 - \left(\frac{1}{2^{n+1}}\right)^{\frac{1}{n}}\right)$ and K decreases in $\left[\frac{1}{c_n}, 1 - \left(\frac{1}{2^{n+1}}\right)^{\frac{1}{n}}\right]$, the inequality $B \leq K(x)$ also holds for $\frac{1}{c_n} \leq x \leq 1 - \left(\frac{1}{2^{n+1}}\right)^{\frac{1}{n}}$. Therefore, we get $1 - \hat{f}_n(x) \leq 1 - f_n(x) + B \leq (1-x)^{c_n} + K(x) = (1-x)^{c_n-1}$. \square

Lemma 6. For $0 \leq x \leq 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$, it holds that $|1 - \hat{f}_n(x)| < (2^n + 1) \cdot \max\{|1-x|^{n+1}, B\}$.

Proof. We first observe that Lemma 4 can be extended from the domain $[0, 1]$ to the larger domain $\left[0, 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1\right]$ when we consider $|1 - f_n(x)|$ and $|1-x|$ instead of $1 - f_n(x)$ and $1-x$ respectively. Assume that $1 < x \leq 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$, and let $H(x) = 2^n \cdot |1-x|^{n+1} - |1-f_n(x)| = 2^n \cdot (x-1)^{n+1} + (-1)^n \cdot (1-f_n(x))$. Then $H'(x) = (n+1)2^n \cdot (x-1)^n - (-1)^n \cdot c_n(1-x^2)^n = (x-1)^n((n+1)2^n - c_n(1+x)^n)$, so there exists a unique local maximal point of H at $x = 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$. Since $H(1) = 0$, it holds that $H(x) \geq 0$ for $1 \leq x \leq 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$. As a result, we obtain

$$|1 - f_n(x)| < 2^n \cdot |1-x|^{n+1}$$

for $0 \leq x \leq 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$. Now we get the result from $|1 - \hat{f}_n(x)| \leq |1 - f_n(x)| + B < 2^n \cdot |1-x|^{n+1} + B \leq (2^n + 1) \cdot \max\{|1-x|^{n+1}, B\}$. \square

Lemma 7. Assume that $B < \frac{1}{2^{n+1}} \cdot \min\left\{\left(\frac{1}{2^{n+1}}\right)^{\frac{1}{n}}, 2\left(\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1\right)\right\}$. If $|1-x| < (2^n + 1) \cdot B$, then it holds that $|1 - \hat{f}_n(x)| < (2^n + 1) \cdot B$.

Proof. Since $B < \frac{2}{2^{n+1}} \cdot \left(\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1\right)$, if $|1-x| < (2^n + 1) \cdot B$, then it holds that $0 < x < 1 + (2^n + 1) \cdot B < 2\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1$. Therefore, we can apply Lemma 6

as following:

$$|1 - \hat{f}_n(x)| < (2^n + 1) \cdot \max \{ (2^n + 1)^{n+1} \cdot B^{n+1}, B \} = (2^n + 1) \cdot B,$$

where the equality comes from $B < \left(\frac{1}{2^n+1}\right)^{\frac{n+1}{n}}$. □

Theorem 6. *Let $B < \frac{1}{2^n+1} \cdot \min \left\{ \left(\frac{1}{2^n+1}\right)^{\frac{1}{n}}, 2 \left(\left(\frac{n+1}{c_n}\right)^{\frac{1}{n}} - 1 \right) \right\}$, and $B < \left(\frac{1}{2^n+1}\right)^{\frac{c_n-1}{n}} - \left(\frac{1}{2^n+1}\right)^{\frac{c_n}{n}}$. For $\epsilon, \alpha > 0$ satisfying $\epsilon \geq \left(\frac{c_n}{c_n-1}\right)^{c_n-1} \cdot B$ and $\alpha \leq \log(1/B) - \log(2^n + 1)$, if $d \geq \frac{1}{\log(c_n-1)} \cdot \log(1/\epsilon) + \frac{1}{\log(n+1)} \cdot \log(\alpha - 1) + O(1)$, then $\|\hat{f}_n^{(d)}(x) - \text{sgn}(x)\|_{[-1, -\epsilon] \cup [\epsilon, 1]} \leq 2^{-\alpha}$.*

Proof. The proof follows the flow of the proof of Theorem 3.

Step 1. It suffices to consider the case $x \in [\epsilon, 1]$ instead of $[-1, -\epsilon] \cup [\epsilon, 1]$. Let $0 < \tau = \left(\frac{1}{2^n+1}\right)^{\frac{1}{n}} < 1$. Our claim is, for any $x \in [\epsilon, 1]$ the inequality $-B < 1 - \hat{f}_n^{(d_\epsilon)}(x) < \tau$ holds for some $0 \leq d_\epsilon \leq d' := \left\lceil \frac{1}{\log(c_n-1)} \cdot \log \left(\log \left(\frac{1}{\tau} \right) / \epsilon \right) \right\rceil$.

Assume that there exists some $x_0 \in [\epsilon, 1]$ that does not satisfy this claim. Since $\epsilon \leq x_0 \leq 1 - \tau$, we obtain $-B < 1 - \hat{f}_n(x_0) < (1 - x_0)^{c_n-1}$ by applying Lemma 5 on x_0 . Then we obtain $\epsilon < x_0 < 1 - (1 - x_0)^{c_n-1} < \hat{f}_n(x_0) < 1 + B < 1 + \tau$. Since $|1 - \hat{f}_n(x_0)| \geq \tau$ by the assumption, it holds that $\epsilon < \hat{f}_n(x_0) \leq 1 - \tau$, so we can apply Lemma 5 on $\hat{f}_n(x_0)$ again which implies $-B < 1 - \hat{f}_n^{(2)}(x_0) < \left(1 - \hat{f}_n(x_0)\right)^{c_n-1}$. By induction, we obtain

$$\begin{aligned} -B < 1 - \hat{f}_n^{(d')} (x_0) &\leq (1 - x_0)^{(c_n-1)d'} \\ &\leq (1 - \epsilon)^{\log(\frac{1}{\tau})/\epsilon} < \left(\frac{1}{\epsilon}\right)^{\log(\frac{1}{\tau})} < \tau, \end{aligned}$$

which contradict to the assumption.

Step 2. Similarly to *Step 1*, we can set our second claim as following: for any $x \in [1 - \tau, 1 + B]$ the inequality $|1 - \hat{f}_n^{(d_\alpha)}(x)| \leq 2^{-\alpha}$ holds for some $0 \leq d_\alpha \leq d'' := \left\lceil \frac{1}{\log(n+1)} \cdot \log \left((\alpha - 1) / \log \left(\frac{1}{(2^n+1)^{\frac{1}{n}} \cdot \tau} \right) \right) \right\rceil$.

Assume that there exists some $x_1 \in [1 - \tau, 1 + B]$ that does not satisfy this claim: $\left|1 - \hat{f}_n^{(d''')} (x_1)\right| \geq 2^{-\alpha} \geq (2^n + 1) \cdot B$ for all $0 \leq d''' \leq d''$. By the assumption, we can say that $x_1 \in [1 - \tau, 1 - (2^n + 1) \cdot B]$, and by applying Lemma 6 on x_1 , we get $|1 - \hat{f}_n(x_1)| \leq (2^n + 1) \cdot (1 - x_1)^{n+1} \leq (2^n + 1) \cdot \tau^{n+1} = \tau$. Therefore, we obtain $1 - \tau \leq \hat{f}_n(x_1) \leq f_n(x_1) + B \leq 1 + B$ so that we can apply

Lemma 6 on $\hat{f}_n(x_1)$. By induction, it holds that

$$\begin{aligned} (2^n + 1)^{\frac{1}{n}} \cdot \left| 1 - \hat{f}_n^{(d'')}(x_1) \right| &\leq \left((2^n + 1)^{\frac{1}{n}} \cdot (1 - x_1) \right)^{(n+1)^{d''}} \\ &\leq \left((2^n + 1)^{\frac{1}{n}} \cdot \tau \right)^{(n+1)^{d''}} \leq 2^{-\alpha+1}, \end{aligned}$$

which contradicts to the assumption.

Combining *Step 1*, *Step 2* and Lemma 7, the proof is completed. □

Corollary 6 (Special Case of Theorem 6 ($n = 4$)). *Let $B < 0.02282$. For $\epsilon \geq 2.15B$ and $\alpha \leq \log(1/B) - 4.09$, if $d \geq 1.83 \log(1/\epsilon) + 0.431 \log(\alpha - 1) + O(1)$, then $\|\hat{f}_4^{(d)}(x) - \text{sgn}(x)\|_{[-1, -\epsilon] \cup [\epsilon, 1]} \leq 2^{-\alpha}$.*

Remark 4. We only addressed about the erroneous evaluation of f_n , but the same logic can be applied to that of g_n : Substituting all c_n 's in Lemma 5 by $g'_n(0)$, then it holds that $-B < 1 - \hat{g}_n(x) < (1 - x)^{g'_n(0)-1}$. As an analogue, by substituting all c_n 's in Theorem 6 by $g'_n(0)$, we can directly convert the theorem into the context of $\hat{f}_n^{(d_f)} \circ \hat{g}_n^{(d_g)}$ instead of $\hat{f}_n^{(d)}$ for $d_g \geq \frac{1}{\log(g'_n(0)-1)} \cdot \log(1/\epsilon) + O(1)$ and $d_f \geq \frac{1}{\log(n+1)} \cdot \log(\alpha - 1) + O(1)$.

One can check that ϵ and α have lower and upper bounds in terms of B , respectively, and this is quite natural: If an input $x > 0$ is so small so that $f_n(x) < B$, then its approximate value $\hat{f}_n(x)$ may be negative due to B -bounded approximation error. Furthermore, if $|x - 1| \ll B$, then $f_n(x)$ should be also very close (even closer) to 1, but a B -bounded approximation error accompanied to $\hat{f}_n(x)$ would disrupt this closeness. In this sense, those lower/upper bounds on ϵ and α with respect to B is actually inevitable.

In fact, Theorem 6 is a *worst-case* analysis on the convergence of $f_n^{(d)}$ in erroneous case by regarding the HEAAN error size in f_n evaluation as B . We also note that inequalities in Lemma 5, 6 and 7 are not as tight as those in Lemma 3 and Lemma 4. In practice, as noted in Sect. 5, even in experiments based on HEAAN, the number of compositions can still be chosen very close to the theoretical lower bounds in Corollary 1 and 3 which are based on the convergence analysis in errorless case.

E Script for Security Estimation

We specified the parameter with security level $\lambda \geq 128$ using the latest LWE estimator [1]⁷. We excluded dec estimates which might not be accurate and often not competitive [21]. The script for checking our parameter is as follows.

⁷ Available on <https://bitbucket.org/malb/lwe-estimator>.

```

load("estimator.py")
n = 2**17; q = 2**3400; alpha = 8/q
duald = partial(drop_and_solve, dual_scale)
primald = partial(drop_and_solve, primal_usvp)
duald(n, alpha, q, secret_distribution=((-1,1), 256),
      reduction_cost_model=BKZ.sieve)
primald(n, alpha, q, secret_distribution=((-1,1), 256),
        rotations=False, reduction_cost_model=BKZ.sieve,
        postprocess=False)

```

References

1. Albrecht, M.R., Player, R., Scott, S.: On the concrete hardness of learning with errors. *J. Math. Cryptol.* **9**(3), 169–203 (2015)
2. Andrievskii, V.: Polynomial approximation of piecewise analytic functions on a compact subset of the real line. *J. Approx. Theory* **161**(2), 634–644 (2009)
3. Armknecht, F., et al.: A guide to fully homomorphic encryption. *Cryptology ePrint Archive*, Report 2015/1192 (2015)
4. Bajard, J.-C., Martins, P., Sousa, L., Zucca, V.: Improving the efficiency of SVM classification with FHE. *IEEE Trans. Inf. Forensics Secur.* **15**, 1709–1722 (2019)
5. Boura, C., Gama, N., Georgieva, M.: Chimera: a unified framework for B/FV, TFHE and HEAAN fully homomorphic encryption and predictions for deep learning. Accepted to *Number-Theoretic Methods in Cryptology (NuTMiC)* (2019)
6. Brakerski, Z.: Fully homomorphic encryption without modulus switching from classical GapSVP. In: Safavi-Naini, R., Canetti, R. (eds.) *CRYPTO 2012*. LNCS, vol. 7417, pp. 868–886. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-32009-5_50
7. Brakerski, Z., Gentry, C., Vaikuntanathan, V.: (Leveled) fully homomorphic encryption without bootstrapping. In: *Proceedings of ITCS*, pp. 309–325. ACM (2012)
8. Chen, H., Chillotti, I., Song, Y.: Improved bootstrapping for approximate homomorphic encryption. In: Ishai, Y., Rijmen, V. (eds.) *EUROCRYPT 2019*. LNCS, vol. 11477, pp. 34–54. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-17656-3_2
9. Cheon, J.H., et al.: Toward a secure drone system: flying with real-time homomorphic authenticated encryption. *IEEE Access* **6**, 24325–24339 (2018)
10. Cheon, J.H., Han, K., Kim, A., Kim, M., Song, Y.: Bootstrapping for approximate homomorphic encryption. In: Nielsen, J.B., Rijmen, V. (eds.) *EUROCRYPT 2018*. LNCS, vol. 10820, pp. 360–384. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-78381-9_14
11. Cheon, J.H., Kim, A., Kim, M., Song, Y.: Homomorphic encryption for arithmetic of approximate numbers. In: Takagi, T., Peyrin, T. (eds.) *ASIACRYPT 2017*. LNCS, vol. 10624, pp. 409–437. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-70694-8_15
12. Cheon, J.H., Kim, D., Kim, D., Lee, H.H., Lee, K.: Numerical method for comparison on homomorphically encrypted numbers. In: Galbraith, S.D., Moriai, S. (eds.) *ASIACRYPT 2019*. LNCS, vol. 11922, pp. 415–445. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-34621-8_15

13. Cheon, J.H., Kim, M., Kim, M.: Search-and-compute on encrypted data. In: Brenner, M., Christin, N., Johnson, B., Rohloff, K. (eds.) FC 2015. LNCS, vol. 8976, pp. 142–159. Springer, Heidelberg (2015). https://doi.org/10.1007/978-3-662-48051-9_11
14. Chialva, D., Dooms, A.: Conditionals in homomorphic encryption and machine learning applications. Cryptology ePrint Archive, Report 2018/1032 (2018)
15. Chillotti, I., Gama, N., Georgieva, M., Izabachène, M.: Faster fully homomorphic encryption: bootstrapping in less than 0.1 seconds. In: Cheon, J.H., Takagi, T. (eds.) ASIACRYPT 2016. LNCS, vol. 10031, pp. 3–33. Springer, Heidelberg (2016). https://doi.org/10.1007/978-3-662-53887-6_1
16. Chillotti, I., Gama, N., Georgieva, M., Izabachène, M.: Faster packed homomorphic operations and efficient circuit bootstrapping for TFHE. In: Takagi, T., Peyrin, T. (eds.) ASIACRYPT 2017. LNCS, vol. 10624, pp. 377–408. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-70694-8_14
17. Comaniciu, D., Meer, P.: Mean shift: a robust approach toward feature space analysis. *IEEE Trans. Pattern Anal. Mach. Intell.* **24**(5), 603–619 (2002)
18. Cordero, A., Soleymani, F., Torregrosa, J.R., Ullah, M.Z.: Numerically stable improved Chebyshev-Halley type schemes for matrix sign function. *J. Comput. Appl. Math.* **318**, 189–198 (2017)
19. Cortes, C., Vapnik, V.: Support-vector networks. *Mach. Learn.* **20**(3), 273–297 (1995)
20. Crawford, J.L., Gentry, C., Halevi, S., Platt, D., Shoup, V.: Doing real work with FHE: the case of logistic regression (2018)
21. Curtis, B.R., Player, R.: On the feasibility and impact of standardising sparse-secret LWE parameter sets for homomorphic encryption. In: Proceedings of the 7th ACM Workshop on Encrypted Computing & Applied Homomorphic Cryptography, pp. 1–10 (2019)
22. 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). https://doi.org/10.1007/978-3-662-46800-5_24
23. Eremenko, A., Yuditskii, P.: Uniform approximation of $\operatorname{sgn} x$ by polynomials and entire functions. *Journal d'Analyse Mathématique* **101**(1), 313–324 (2007)
24. Fan, J., Vercauteren, F.: Somewhat practical fully homomorphic encryption. *IACR Cryptology ePrint Archive*, 2012:144 (2012)
25. Friedman, J.H.: Greedy function approximation: a gradient boosting machine. *Ann. Stat.* **29**, 1189–1232 (2001)
26. Friedman, J.H.: Stochastic gradient boosting. *Comput. Stat. Data Anal.* **38**(4), 367–378 (2002)
27. Gentry, C.: A fully homomorphic encryption scheme. Ph.D. thesis, Stanford University (2009). <http://crypto.stanford.edu/craig>
28. 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. LNCS, vol. 8042, pp. 75–92. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-40041-4_5
29. Gilad-Bachrach, R., Dowlin, N., Laine, K., Lauter, K., Naehrig, M., Wernsing, J.: Cryptonets: applying neural networks to encrypted data with high throughput and accuracy. In: International Conference on Machine Learning (2016)
30. Goldschmidt, R.E.: Applications of division by convergence. Ph.D. thesis, Massachusetts Institute of Technology (1964)

31. Han, K., Hong, S., Cheon, J.H., Park, D.: Logistic regression on homomorphic encrypted data at scale. In: The AAAI Conference on Innovative Applications of Artificial Intelligence (2019)
32. Han, K., Ki, D.: Better bootstrapping for approximate homomorphic encryption. Cryptology ePrint Archive, Report 2019/688 (2019). To Appear in CT-RSA 2020
33. Hartigan, J.A., Wong, M.A.: Algorithm as 136: a k-means clustering algorithm. *J. Royal Stat. Soc. Ser. C (Appl. Stat.)* **28**(1), 100–108 (1979)
34. Higham, N.J.: *Functions of matrices: theory and computation*. SIAM (2008)
35. Kazarinoff, D.K.: On Wallis' formula. *Edinb. Math. Notes* **40**, 19–21 (1956)
36. Kenney, C.S., Laub, A.J.: The matrix sign function. *IEEE Trans. Autom. Control* **40**(8), 1330–1348 (1995)
37. Kim, A., Song, Y., Kim, M., Lee, K., Cheon, J.H.: Logistic regression model training based on the approximate homomorphic encryption. *BMC Med. Genomics* **11**(4), 83 (2018)
38. Kim, D., Son, Y., Kim, D., Kim, A., Hong, S., Cheon, J.H.: Privacy-preserving approximate GWAS computation based on homomorphic encryption. Cryptology ePrint Archive, Report 2019/152 (2019)
39. Kim, M., Song, Y., Li, B., Micciancio, D.: Semi-parallel logistic regression for GWAS on encrypted data. Cryptology ePrint Archive, Report 2019/294 (2019)
40. Lin, Y.: A note on margin-based loss functions in classification. *Stat. Probab. Lett.* **68**(1), 73–82 (2004)
41. Mitrinović, D.S., Pečarić, J.E., Fink, A.: Bernoulli's inequality. In: Mitrinović, D.S., Pečarić, J.E., Fink, A. (eds.) *Classical and New Inequalities in Analysis*, pp. 65–81. Springer, Dordrecht (1993). <https://doi.org/10.1007/978-94-017-1043-5>
42. Nakatsukasa, Y., Bai, Z., Gygi, F.: Optimizing Halley's iteration for computing the matrix polar decomposition. *SIAM J. Matrix Anal. Appl.* **31**(5), 2700–2720 (2010)
43. Paterson, M.S., Stockmeyer, L.J.: On the number of nonscalar multiplications necessary to evaluate polynomials. *SIAM J. Comput.* **2**(1), 60–66 (1973)
44. Saff, E., Totik, V.: Polynomial approximation of piecewise analytic functions. *J. London Math. Soc.* **2**(3), 487–498 (1989)
45. Snucrypto. HEAAN (2017). <https://github.com/snucrypto/HEAAN>
46. Soheili, A.R., Toutounian, F., Soleymani, F.: A fast convergent numerical method for matrix sign function with application in SDEs. *J. Comput. Appl. Math.* **282**, 167–178 (2015)
47. Tan, B.H.M., Lee, H.T., Wang, H., Ren, S.Q., Khin, A.M.M.: Efficient private comparison queries over encrypted databases using fully homomorphic encryption with finite fields. *IEEE Trans. Dependable Secure Comput.* (2020)
48. Wilkes, M.V.: *The Preparation of Programs for an Electronic Digital Computer: With special reference to the EDSAC and the Use of a Library of Subroutines*. Addison-Wesley Press (1951)