# FIDESlib: A Fully-Fledged Open-Source FHE Library for Efficient CKKS on GPUs

Carlos Agulló-Domingo
*Universidad de Murcia*
carlos.a.d@um.es

Óscar Vera-López
*Universidad de Murcia*
oscar.veral@um.es

Seyda Guzelhan
*Boston University*
seyda@bu.edu

Lohit Daksha
*Boston University*
tihol@bu.edu

Aymane El Jerari
*Northeastern University*
eljerari.a@northeastern.edu

Kaustubh Shivdikar
*AMD*
kshivdik@amd.com

Rashmi Agrawal
*Boston University*
rashmi23@bu.edu

David Kaeli
*Northeastern University*
kaeli@ece.neu.edu

Ajay Joshi
*Boston University*
joshi@bu.edu

José L. Abellán
*Universidad de Murcia*
jlabellan@um.es

*Abstract*—Word-wise Fully Homomorphic Encryption (FHE) schemes, such as CKKS, are gaining significant traction due to their ability to provide post-quantum-resistant, privacy-preserving approximate computing—an especially desirable feature in Machine-Learning-as-a-Service (MLaaS) cloud-computing paradigms. OpenFHE is a leading CPU-based FHE library with robust CKKS operations, but its server-side performance is not yet sufficient for practical cloud deployment. As GPU computing becomes more common in data centers, many FHE libraries are adding GPU support. However, integrating an efficient GPU backend into OpenFHE is challenging. While OpenFHE uses a Hardware Abstraction Layer (HAL), its flexible architecture sacrifices performance due to the abstraction layers required for multi-scheme and multi-backend compatibility.

In this work, we introduce FIDESlib, the first open-source server-side CKKS GPU library that is fully interoperable with well-established client-side OpenFHE operations. *Unlike other existing open-source GPU libraries*, FIDESlib provides the first implementation featuring heavily optimized GPU kernels for all CKKS primitives, including bootstrapping. Our library also integrates robust benchmarking and testing, ensuring it remains adaptable to further optimization. Furthermore, its software architecture is designed to support extensions to a multi-GPU backend for enhanced acceleration. Our experiments across various GPU systems and the leading open-source CKKS library to date, Phantom, show that FIDESlib offers superior performance and scalability. For bootstrapping, FIDESlib achieves no less than $70\times$ speedup over the AVX-optimized OpenFHE implementation.

## I. INTRODUCTION AND MOTIVATION

Machine-Learning-as-a-Service (MLaaS) [1] is a cloud-based powerful tool to simplify the development, deployment, and management of machine learning models. Nowadays, major cloud providers enable robust ML platforms such as Amazon SageMaker [2] for data preparation, training, and deployment; Google Vertex AI [3] for AutoML and custom training; and Azure Machine Learning [4] for automated end-to-end workflows. While MLaaS platforms are gaining traction, data breaches pose significant risks, leading to millions in losses [5]. This threat is further amplified by the anticipated rise of quantum computers, which are expected to possess encryption-breaking capabilities within the next decade [6]. Fortunately, lattice-based cryptography, endorsed by NIST,

offers a solution, with post-quantum-resistant standards such as FIPS 203 [7] and FIPS 204 [8] recently announced [9].

CKKS [10] is a modern, Ring-Learning-With-Errors (RLWE) lattice-based cryptographic scheme that is emerging as the preferred Fully Homomorphic Encryption (FHE) technique for privacy-preserving approximate computing (PPaC), thereby potentially addressing the critical data breach challenges in MLaaS-enabled cloud platforms. However, the large computational overhead of CKKS incurs a 2–5 orders of magnitude slowdown [11], hindering its practical deployment in many of these cloud services, such as real-time ML inference.

OpenFHE [10] is the leading CPU-based FHE library, extensively supported by a large, multidisciplinary research community. As a result, OpenFHE incorporates robust and rigorously validated client-side CKKS operations (e.g., encryption/descryption) and server-side CKKS operations (e.g., modular arithmetic and bootstrapping). However, despite its AVX-optimized CPU backend [12], the server-side CKKS operations remain far from achieving the level of performance required for practical deployment in this class of cloud services. Although FPGA and ASIC CKKS accelerators [13] could achieve the desired practicality, their widespread adoption as established compute nodes in data centers faces significant challenges. These include the high integration costs within existing and mature MLaaS ecosystems, making it difficult for them to become a reality.

Given the growing presence of GPU-based computing nodes in high-performance data centers [14], accelerating modern CKKS-powered cloud services by fully leveraging the capabilities of GPU platforms is becoming increasingly common. Consequently, most state-of-the-art CKKS libraries are already utilizing GPUs [15]–[23]. However, while adding a GPU backend to OpenFHE would mark significant progress, its multi-scheme and multi-backend general software architecture inherently sacrifices cutting-edge performance in favor of generality, due to the many abstraction layers needed for such flexibility.

State-of-the-art GPU libraries that efficiently support all CKKS functionality are often kept private for internal use or industry applications [16], [17], [22]. This lack of open ac-

cess limits collaboration across the broader multi-disciplinary research community, hindering the development of new algorithmic optimizations aligned with the rapid advancements in GPU architectures. It also restricts the adoption of CKKS in MLaaS scenarios and prevents general users from benefiting from these optimizations with privacy-preserving computation guarantees. To address these challenges, democratizing CKKS optimization research for modern GPUs is essential.

To achieve this important goal, we propose FIDESlib, the first open-source CKKS GPU library that is *feature-complete*, including bootstrapping, and overcomes the limitations of existing open-source counterparts, delivering the highest performance. To facilitate collaborative development and further optimization of our library, our design strategy for FIDESlib is threefold. First, we leverage the well-established and robust OpenFHE library for all client-side operations. Second, we implement all CKKS server-side operations, adding exhaustive testing capabilities (unit and integration tests) for complete functional validation against OpenFHE's server operations, using the GoogleTest framework [24]. Third, we enable extensive microbenchmarking and benchmarking capability by utilizing the Google Benchmark framework [25].

Our key contributions are as follows.

- We propose FIDESlib[1], the first open-source implementation of CKKS with full functionality, including bootstrapping, specifically optimized to maximize compute throughput, on-chip data reuse, and memory bandwidth utilization for a GPU backend.
- FIDESlib is the first GPU library that provides full interoperability with the standard OpenFHE framework, enabling the implementation of server-side CKKS operations while delegating encoding and encryption tasks to OpenFHE-based clients.
- Our experiments across various GPU platforms show that FIDESlib outperforms Phantom [15], the leading open-source CKKS library, in performance and scalability for all operations supported by Phantom. For CKKS operations not supported by Phantom, such as bootstrapping, FIDESlib achieves up to $74\times$ speedup over the best-performing AVX-optimized OpenFHE implementation.

## II. BACKGROUND

### A. CKKS Optimization Techniques

HE enables operations on encrypted data without having access to the secret key. Following Gentry's blueprint [26], several FHE schemes have been developed to support exact arithmetic on integers (BGV [27] and BFV [28]), approximate arithmetic on fixed-point numbers (CKKS [10]), and Boolean gates (TFHE [29] and FHEW [30]). Given CKKS's potential to enable PPaC in MLaaS cloud environments, we have implemented and optimized this scheme in FIDESlib, targeting a GPU backend. Table I lists the available primitives of the CKKS cryptosystem, while Table II presents the main CKKS

| Operation | Output | Description |
|---|---|---|
| ScalarAdd($\|x\|, c$) | $\|x+c\|$ | Adding a constant to a ciphertext. |
| PtAdd($\|x\|, y$) | $\|x+y\|$ | Addition of a plaintext with a ciphertext. |
| HAdd($\|x\|, \|y\|$) | $\|x+y\|$ | Addition of two ciphertexts. |
| ScalarMult($\|x\|, c$) | $\|x \odot c\|$ | Multiplying a ciphertext by a constant. |
| PtMult($\|x\|, y$) | $\|x \odot y\|$ | Multiplying a ciphertext by a plaintext. |
| HMult($\|x\|, \|y\|$) | $\|x \odot y\|$ | Multiplication of two ciphertexts. |
| Rescale($\|x\|_{C_L}$) | $\|x\|_{C_{L-1}}$ | Rescale after multiplication. |
| Conjugate($\|x\|$) | $\|\bar{x}\|$ | Conjugation of the underlying message. |
| HRotate($\|x\|, k$) | $\|\phi_k(x)\|$ | Rotation of the underlying message by k. |

TABLE I: CKKS RNS primitives.

| Param | Description |
|---|---|
| $N$ | Polynomial degree-bound |
| $n$ | Length of the message. $n \le \frac{N}{2}$ |
| $Q$ | Polynomial modulus |
| $L$ | Multiplicative depth before bootstrapping is needed |
| $\mathcal{B}$ | The set $\{q_0, q_1, \ldots, q_L\}$ of prime factors of $Q$ |
| $\ell + 1$ | Current number of limbs |
| $\ell$ | remaining multiplicative depth |
| dnum | Number of digits in the switching key |
| $P$ | Product of extension limbs added for raised modulus. |
| $\Delta$ | scaling applied during encoding. $q_i \approx \Delta$ |
| $x$ | a polynomial |
| $[x]_{q_i}$ | $q_i$-limb of $x$ |
| **ksk** | Evaluation/key-switching key |
| NTT | Number Theoretic Transform, takes a coefficient vector and outputs an evaluation vector |
| iNTT | Inverse Number Theoretic Transform, takes an evaluation vector and outputs a coefficient vector |
| $\psi$ | $2n$-th root of unity under a prime $q_i$, typically precomputed for (i)NTT evaluation. |

TABLE II: CKKS Notation and descriptions.

parameters and their notation used in this work. For a detailed description of the CKKS scheme, refer to [31].

In contrast to BFV and BGV, CKKS treats the fresh encryption error as part of the noise from fixed-point operations [10]. The most well-known CKKS libraries (HElib [32], HeaNN [10], SEAL [33], OpenFHE [34]) provide efficient implementations of FHE on CPU platforms. Several algorithmic optimizations exist to reduce computation costs and enhance CKKS precision, which can also be applied to a GPU backend. Below, we discuss the most significant of these optimizations.

The large size of the ciphertext modulus can be resolved by using the Residue Number System (RNS) technique based on the Chinese Remainder Theorem [35]. Precision of rescaling under RNS can be improved by carefully tracking the scaling factors at each level [36].

Key-Switching is needed by HMult and HRotate operations to keep the encrypting key of the result the same as the input by employing a "switching key". Hybrid key switching proposes the idea of decomposing the switching keys into $dnum$ chunks so that the noise handling is easier in each small partition [37]. This comes with a trade-off of increased computation and key-size.

The number of multiplicative operations on a ciphertext is limited by the amount of accumulated noise. The bootstrapping

operation is used to refresh the ciphertext [38] noise and allow for further computation. Algorithmic improvements to bootstrapping focus mainly on linear transforms [39], [40] or polynomial evaluation [41], or both [37], [42], [43]. The linear transforms can be improved by exploiting matrix sparsity [44] and applying decomposition methods [40] to the original DFT matrix. Even with these optimizations, bootstrapping consumes more than half of the runtime of a typical FHE-enabled application when run on a GPU [19].

### B. The Memory Bottleneck and GPUs

As explored in many prior studies [45]–[47], FHE workloads are heavily memory-bound. This is primarily due to the combination of low arithmetic intensity operations, large ciphertext and key size, and lack of advances in memory bandwidth and latency in modern compute platforms [48]. As a result, FHE workloads require new acceleration strategies. Computations performed in the ciphertext domain are typically 2 to 5 orders of magnitude slower than their plaintext counterparts [11]. This degree of slowdown is due to the use of large polynomials that encode plaintext data used in ciphertext computations. Large polynomial parameters need to be used to maintain adequate levels of security, as well as to ensure noise growth does not impact the underlying data. For example, to achieve 128-bit security, FHE operations require using polynomials of a degree up to $2^{17}$, with the size of each coefficient being up to 2200 bits [49].

Unlike CPU-based FHE implementations, which are constrained by the bandwidth of DDR memory, GPUs can leverage GDDR or High Bandwidth Memory (HBM) modules, helping to alleviate the memory bottleneck. While modern DDR5-based systems can achieve a peak memory bandwidth of around 100GB/s, consumer GPUs like the RTX 4090 used in our experiments exceed 1 TB/s of memory bandwidth, significantly accelerating CKKS primitives. Additionally, the on-chip memory capacity in GPUs plays a crucial role in data reuse effectiveness. When properly utilized, it reduces costly off-chip memory accesses. For instance, the RTX 4090 includes a large 72MB L2 cache, which will be essential for further acceleration in FIDESlib. Although CKKS primitives are not typically compute-intensive, their performance on GPUs can be significantly affected by modular arithmetic operations on integers (the key operations in CKKS), as current GPU architectures lack efficient hardware support for modulo operations. To demonstrate the efficiency of our FIDESlib library across different GPU platform capabilities, we will test four distinct GPU platforms in Section IV (listed in Table IV).

## III. FIDESLIB

### A. Functionality Overview

FIDESlib implements all server-side CKKS operations (e.g., HAdd, HMult) on the GPU, while client-side operations (e.g., KeyGen, Encrypt, Decrypt) are handled by the OpenFHE library. Figure 1 shows how functionality is divided into the two libraries. To the best of our knowledge, FIDESlib includes

the first open-source GPU implementation of the CKKS bootstrapping procedure. FIDESlib also offers optimized routines such as ScalarAdd, ScalarMult and HSquare, which are more efficient versions of PtAdd, PtMult, and HMult for inputs with repetitive data (features rarely available in other libraries). Additionally, the HoistedRotation routine [39] is implemented to reduce the cost of multiple ciphertext rotations.
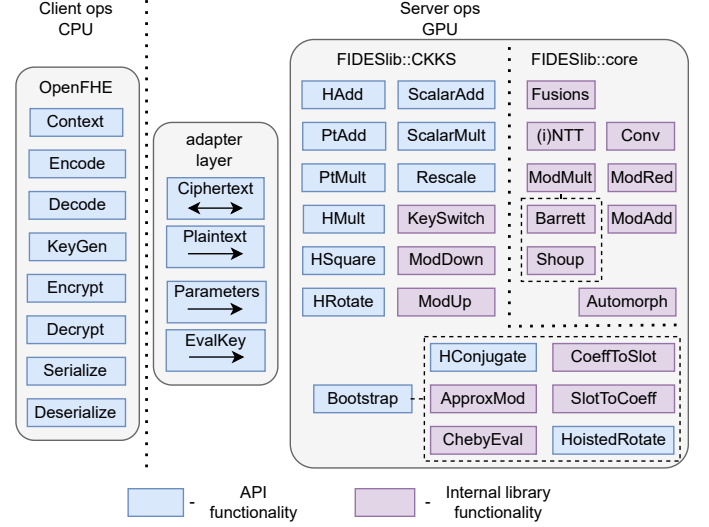


Fig. 1: High-Level Functionality Overview in FIDESlib.

### B. OpenFHE interoperation

OpenFHE adheres to the *HomomorphicEncryption.org* post-quantum security standards for homomorphic encryption. As the security of the library depends solely on the correctness of client-side operations, users of the FIDESlib library benefit from the same security guarantees. However, the OpenFHE and FIDESlib source codes are decoupled by implementing a thin adapter layer. Loose coupling offers several benefits:

- **Changes to OpenFHE**: Only minimal changes are made to the OpenFHE codebase, ensuring that new versions of the OpenFHE library are unlikely to break compatibility. In our repository, the instructions for compiling a compatible version of OpenFHE for FIDESlib include applying a small Git patch to any OpenFHE release.
- **Freedom and Simplicity of Implementation**: The generic multi-backend design of OpenFHE's implementation of multi-FHE schemes makes developing a performant GPU acceleration backend challenging. By adopting our own code structure and focusing on the CKKS scheme, we reduce code complexity while maintaining the flexibility required for further optimization.
- **Code Repurposing**: Since the core functionality of the library is not tied to OpenFHE's code, FIDESlib can be easily adapted as a GPU backend for other projects.

Internally, the adapter layer transfers data between OpenFHE's objects and simplified data structures that retain essential data and metadata fields. These simplified data structures are then passed to methods within FIDESlib's
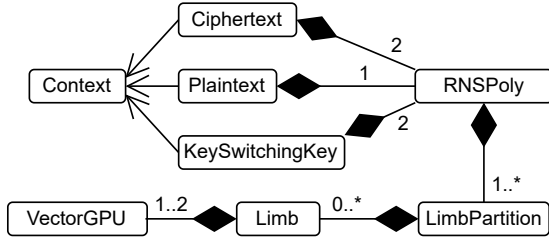
Fig. 2: Data class diagram.

classes, such as *CKKS::Parameters*, *CKKS::Ciphertext*, *CKKS::Plaintext*, and *CKKS::KeySwitchingKey*, which store the data in GPU memory. For the *CKKS::Ciphertext* class, the reverse process is also supported, transferring the data, along with static noise estimation data, back to an OpenFHE ciphertext object for decryption.

### C. Software design

The FIDESlib codebase is organized into four namespaces:

1) **FIDESlib**: core functionality and utility code to support power of two polynomial ring arithmetic under word-sized moduli on GPU.
2) **FIDESlib::CKKS**: implementation of the CKKS's functionality, including GPU kernels, a composition-based hierarchy of GPU-memory-managing classes, and the CKKS crypto-context class.
3) **FIDESlib::test**: a Google Test-based test suite comprising both unit and integration tests.
4) **FIDESlib::bench**: a Google Bench-based benchmark suite for testing the library's functionality and performance with synthetic and real workloads.

In the remainder of this section, we focus on the **FIDESlib** and **FIDESlib::CKKS** namespaces. For brevity, the **FIDESlib::CKKS** prefix will be omitted.

The organization of data within the library is illustrated in Figure 2, which is discussed below.

*Ciphertext*, *Plaintext*, and *KeySwitchingKey* are composed of a number of *RNSPoly* objects. Each *RNSPoly* represents an $N$ degree polynomial, as determined by its generating *Context*. The polynomial is defined under a specific integer modulus and decomposed into a base $\mathcal{B} = q_0, q_1, ..., q_l$, as specified by the Remainder Number System.

Each *RNSPoly* object is composed of one or more *LimbPartition* objects, which represent the portion of a polynomial stored in a specific device. As the current FIDESlib version supports single-GPU operations (with multi-GPU support in development), each *RNSPoly* object can be assumed to contain one *LimbPartition*, with all *LimbPartition* instances assigned to the same GPU.

Each *LimbPartition* contains one or several *stacks* of *Limb* objects, which represent a polynomial under a specific base modulus $q_i \in \mathcal{B}$. The *LimbPartition* class accesses metadata arrays specific to each *Limb stack* type it stores and provides methods for managing them. These metadata stacks are defined

during *Context* creation to simplify precomputation and GPU execution setup.

Finally, a *Limb* instance contains a *VectorGPU* object that manages a contiguous array of memory and an optional second *VectorGPU* object for $(i)NTT$ computation. The *Limb* class is templated to support both 64-bit and 32-bit integers, although the current version only verifies and utilizes 64-bit integers.

### D. Memory management

FIDESlib employs a stack of arrays approach instead of using a flattened 2D array of coefficients to represent an RNS polynomial. The main benefit of this approach is its finer grain memory management, which can reduce internal fragmentation. Additionally, we argue that it simplifies programming.

The *VectorGPU* class simplifies memory management using the **RAII** (Resource Acquisition Is Initialization) pattern, asynchronously allocating device memory at construction and asynchronously freeing it at destruction—i.e., when the object goes out of scope—removing the need for explicit device memory management by the programmer. Asynchronous device memory management is achieved through the CUDA Stream Ordered Memory Allocator, a.k.a. the default *cudaMemPool_t*.

Although, the stack-of-arrays representation simplifies usage for library programmers, we also offer the option to create unmanaged *VectorGPU* objects. These are instantiated by explicitly passing the device data pointer to the constructor, shifting memory management responsibility to a higher-level class. We observe that for short-lived or constant-sized *RNSPoly* objects, creating a flattened 2D data array while simulating a stack of arrays through the use of unmanaged *VectorGPU*s can be preferable. In contrast, long-lived and variably-sized *RNSPoly* objects are more easily and memory efficiently managed as true stacks of arrays. Thus, FIDESlib employs both approaches.

### E. Precomputation

Many values are precomputed for the CKKS scheme before server operations are executed, reducing the computational complexity of many operations. For simplicity, most of this precomputation is done during *Context* creation time. Next, we discuss these precomputed values and their implications for algorithm implementation.

In OpenFHE, multiple cryptocontexts can be created simultaneously. However, two limitations of the GPU platform make this impractical. The first limitation is the limited size of device memory. The second is that constant memory fields must be defined at compile-time and are restricted to an aggregated size of 64 KB. Running a scheme with different parameter sets concurrently becomes challenging to implement and impractical for a GPU backend.

Consequently, FIDESlib treats Contexts and precomputation results using the singleton software design pattern. This approach allows for the use of constant memory and enables precomputation values to be declared globally. While this may seem unorthodox, it simplifies kernel calls, as the precomputed values do not need to be passed as kernel parameters. Note

that application developers can bypass this library limitation by managing different parameter sets within separate processes, where FIDESlib is linked statically. However, this approach does not resolve the issue of device memory usage.

Finally, values to be used by every thread in a *warp* simultaneously are declared *__constant__*, as the hardware efficiently broadcasts this value to all threads. For values that are worth precomputing but not necessarily accessed by all threads at once, they are declared *__global__*.

### F. Algorithms and optimizations

In this subsection, we describe the core algorithms and optimizations implemented in FIDESlib.

*1) Stream Handling and Limb Batching:* Computational kernels in the CKKS scheme can be categorized based on their data dependencies:

- **No dependencies**: These include elementwise operations such as modular multiplication, modular addition, and modulus switching.
- **Dependencies amongst elements within a limb**: Operations like NTT, iNTT, and automorphism are dependent on the data within a single limb.
- **Dependencies amongst elements in the same index of different limbs**: Operations such as ModUp, ModDown, and Rescale involve dependencies across elements in different limbs at the same index.

Given these dependencies, the minimum working set size for complex CKKS operations, such as HMult, corresponds to an entire ciphertext plus a key switching key, which is typically very large (e.g. 120 MB) and challenging to retain on-chip. However, the GPU's L2 cache hit rate can be improved by performing consecutive operations on a subset of a ciphertext's limbs before the data is evicted to main memory. This approach has been shown to deliver significant performance gains [45]. Nonetheless, working with only a subset of a ciphertext's limbs may lead to under-utilization of the GPU's compute throughput.

To tackle this problem, FIDESlib divides operations without dependencies across different limbs into multiple kernels, with each kernel processing one or more limbs independently. These independent kernels run asynchronously in separate CUDA streams, maximizing compute utilization. Using kernels with smaller memory footprints enhances temporal locality in data access, thereby improving L2 cache hit rates.

While assigning a single limb per kernel improves locality, the CPU's kernel launch overhead becomes a bottleneck on faster GPUs, which complete tasks quicker than the CPU can schedule. To mitigate this, FIDESlib introduces "limb-batching" controlled by a configurable *Context* parameter, to reduce CPU overhead.

*2) Modular arithmetic:* Since integer operations on limb elements must be performed modulo a specific set of prime numbers $\{p_i\}_0^{L+K}$ selected at runtime, it is essential to optimize the modulo operation because it is not natively supported by GPUs. Note that using the modulo operator % naively in CUDA code leads to complex compiler generated

| Method | Multiplications | Output | Observation |
|---|---|---|---|
| Montgomery reduction | 1 (wide) + 1 (low) | $[0, 2p)$ | Requires inputs to be in Montgomery format |
| Montgomery multiplication | 2 (wide) + 1 (low) | | |
| Shoup multiplication | 1 (wide) + 2 (low) | $[0, 2p)$ | Precomputation depends on one of the inputs |
| Barrett reduction | 1 (wide) + 1 (low) | $[0, 2p)$ | |
| Barrett multiplication | 2 (wide) + 1 (low) | | |

TABLE III: Comparison of fast modular reduction methods. Montgomery, Shoup, and improved Barrett [50]. Wide multiplications (64-bit $\times$ 64-bit $\rightarrow$ 128-bit) are considerably more expensive than low multiplications (64-bit $\times$ 64-bit $\rightarrow$ 64-bit).

routines, which require numerous assembly instructions when the modulo operand is not known at compile time.

Modular addition and subtraction are relatively simple to implement, as their results fall within the ranges $[0, 2p-1)$ and $(-p, p)$, respectively, and can be easily adjusted back to $[0, p)$. In contrast, reducing the results of multiplication operations is significantly more challenging, as they fall within $[0, p^2 - p)$.

Several optimized modular reduction routines have been developed, and Table III compares some of the most commonly used techniques in the literature. Barrett and Montgomery reduction methods improve the efficiency of modular reduction by precomputing auxiliary values dependent on the modulus. The Shoup technique, while more computationally efficient, requires precomputing a constant based on both the modulus and one of the operands. FIDESlib implements the improved Barrett reduction [50], which offers the same computational efficiency as the Montgomery technique but does not require a specific encoding of the elements. Additionally, FIDESlib leverages Shoup's modular multiplication where appropriate, further accelerating certain modulo operations.

*3) The base conversion kernel:* The main algorithmic component of the ModDown, ModUp, and Rescale operations is the fast base conversion algorithm. Rescale is a special case where the CRT base is modified by a single prime number, so it is optimized separately. ModDown and ModUp, on the other hand, are used in the (Hybrid) KeySwitching procedure during HMult and HRotate.

Computationally, the fast base conversion algorithm $(Conv_{\mathcal{B}' \rightarrow \mathcal{B}}([x(X)]_{\mathcal{B}'}) = [x(X)]_{\mathcal{B}})$ can be thought of as a modular matrix-vector multiplication, preceded by a limb-wise scaling of the coefficients, as shown by Equation 1. The same computation is applied to every coefficient of the limb, making it a matrix-matrix product.

$$\begin{pmatrix} x(X)^{(0)} \\ x(X)^{(1)} \\ \vdots \\ x(X)^{(K-1)} \end{pmatrix} = \begin{pmatrix} [\hat{q}_0]_{P_0} & \cdots & [\hat{q}_{L-1}]_{P_0} \\ [\hat{q}_0]_{P_1} & \cdots & [\hat{q}_{L-1}]_{P_1} \\ \vdots & \ddots & \vdots \\ [\hat{q}_0]_{P_{K-1}} & \cdots & [\hat{q}_{L-1}]_{P_{K-1}} \end{pmatrix} \cdot \begin{pmatrix} \hat{q}_0^{-1} \cdot x(X)^{(0)} \\ \hat{q}_1^{-1} \cdot x(X)^{(1)} \\ \vdots \\ \hat{q}_{L-1}^{-1} \cdot x(X)^{(L-1)} \end{pmatrix}$$

(1)

As digit decomposition is applied before base conversion in ModUp, the original base, with $(L+1)/dnum$ elements, becomes up to $dnum$ times smaller than the output base, which has $l+1$ elements. Typically, values for $L/dnum$ are small, so we cache the result of the initial scaling applied to a subset of the limb coefficients in shared memory. We then accumulate each dot product and write the results back to main memory. Two GPU threads are used per vector-matrix operation, maximizing compute resource utilization and ensuring optimal performance.

The kernel thus utilizes $4 \cdot L/dnum$ shared memory bytes per thread. Each shared memory element being reused $l+1$ times, the number of remaining limbs that remain as multiplicative levels are used.

To optimize performance, partial dot product results are accumulated as a 128-bit integer and are only reduced modulo $p_i$ before being written back to global memory, saving $l$ reductions per output element. Nonetheless, the base conversion kernel remains compute-bound, especially with large parameter sets.

*4) (i)NTT:* The (i)NTT kernel enables efficient transformation of polynomials between the coefficient and evaluation spaces. Additions can be performed in either representation, while polynomial multiplication is efficiently carried out in the evaluation space. Other transformations, such as base conversion, must be executed in the coefficient representation.

Internally, the NTT algorithm is implemented as a nega-cyclic convolution of polynomial coefficients, similar to a Discrete Fourier Transform (DFT) calculation but performed over an integer ring rather than complex numbers. Consequently, the NTT can be implemented using a modified version of the FFT algorithm. Given that the ring dimension $N$ is defined as a power of two, FIDESlib employs the Radix-2 FFT scheme. While larger radices have been proposed (Radix-8 [15]), the Radix-2 algorithm minimizes computational complexity, which we found to be the primary bottleneck.

As $N \in 2^{13,\ldots,17}$ for typical parameter sets, a single limb $(64-512\text{KB})$ does not fit into the shared memory of a Streaming Multiprocessor (SM). Therefore, the algorithm must rely on global memory accesses for communication. To minimize these accesses, FIDESlib employs a Hierarchical/2D NTT scheme. This approach reduces the NTT computation to just four memory accesses per element, significantly fewer than a naive 1D NTT implementation. Figure 3 provides an overview of the data movement involved in our implementation. As accesses to the coefficient vector must be performed with a stride, each thread block performs computation on four $\sqrt{N}$-sized sub-FFTs. This strategy leverages the memory block size of the GPU architecture. Additionally, performing multiple FFTs within a single thread block helps amortize the cost of loading the precomputed twiddle values $(\psi_i)$ into shared memory.

To reduce L2 load bandwidth usage, the required full vector scaling steps compute twiddle factors "on the fly" using some of the input from the already loaded twiddle factor vector. By precomputing the twiddle factors, corresponding Shoup
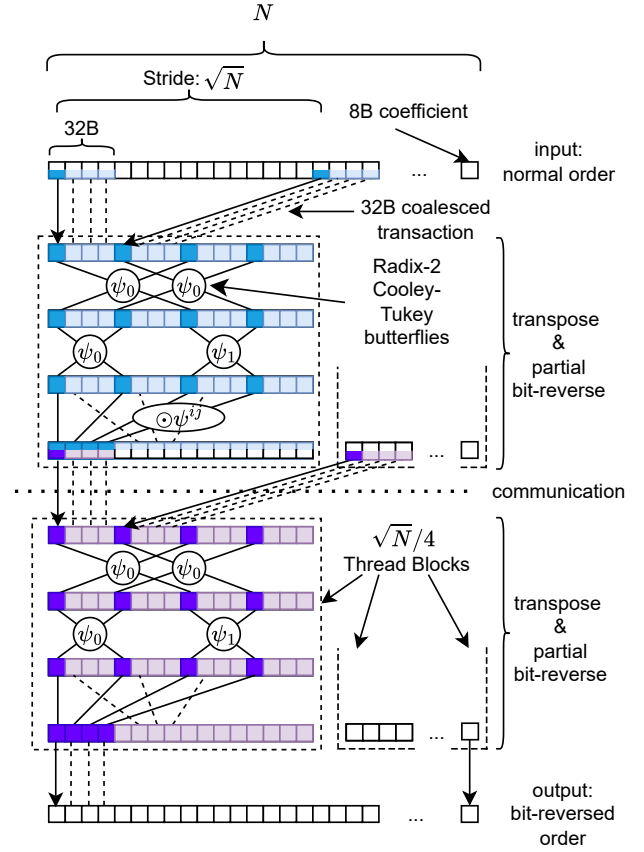


Fig. 3: Data movement diagram of the Radix-2 2D/Hierarchical NTT algorithm implemented in FIDESlib.

constants are also precomputed, enabling Shoup modular multiplication to speed up the NTT butterflies.

The iNTT is implemented with Gentleman-Sande butterflies instead of Cooley-Tukey butterflies, which take a bit-reversed evaluation polynomial and output a normal-ordered coefficient polynomial. This approach eliminates the need for explicit bit-reversal steps.

*5) Operation Fusions:* A common technique employed to reduce the memory bandwidth requirements of CKKS is kernel fusion. As the name suggests, several mathematical steps are computed within a single kernel, which reduces the frequency of data reads from and writes to global memory.

Before showcasing some instances of operation fusion employed by FIDESlib, we will discuss the implementation. The main drawback of implementing operation fusion is the combinatorial growth of possible kernel fusion combinations. The primary kernel fusion approach involves appending element-wise operations to the more compute-intensive (i)NTT kernels. To facilitate this, our (i)NTT kernels use a template to include kernel fusions without duplicating any code.

- **Rescale fusion**: The NTT kernel starts with a fused $SwitchModulo$ operation and ends with the combined sequence: $q_l^{-1}(x^{(i)} - NTT(SwitchModulo(x^{(l)})))$, avoiding up to six memory operations per element.

- **ModDown fusion**: Similarly, the NTT kernel performs the $\mathcal{P}^{-1}(x^{(i)} - NTT(x'^{(i)}))$ sequence, saving up to four memory operations per element.
- **HMult fusion**: The outputs $\{\mathcal{P} * (c_0^{(i)} \odot c_0'^{(i)}) * ksk_{0,d_i}^{(i)}, \mathcal{P} * (c_0^{(i)} \odot c_1'^{(i)} + c_1^{(i)} \odot c_0'^{(i)}) \odot ksk_{1,d_i}^{(i)}, INTT(c_1^{(i)} \odot c_1'^{(i)})\}$ are produced by the iNTT kernels before ModUp saving up to 10 memory operations per element. After the base conversion step, the NTT kernels generate $\{NTT(x'^{(i)}) * ksk_{0,\beta}^{(i)}, NTT(x'^{(i)}) * ksk_{1,\beta}^{(i)}\}$ as a result, saving up to three memory operations per element.
- **Dot product fusion**: The pattern $\sum_{i=0}^{n-1} a \odot b$ can be optimized from $6n - 3$ memory operations per output element to $2n + 1$ by computing everything in a single kernel. Similarly, a weighed sum can be reduced from $4n - 2$ down to $n + 1$ memory operations per element.

*6) Hoisted rotation:* When multiple rotations of a single ciphertext are required, the ModUp operation is performed once before applying key switching with different keys to its output [39]. This optimization reduces the overall runtime to obtain the rotated ciphertexts. FIDESlib implements this routine, which is especially useful for accelerating Bootstrap execution.

*7) Bootstrapping:* FIDESlib adapts the CKKS bootstrapping procedure from OpenFHE for execution on GPU. The bootstrapping runtime consists of three main routines: CoeffToSlot, ApproxModEval, and SlotToCoeff [38].

FIDESlib adapts the ApproxModEval routine directly from OpenFHE [43]. It implements Chebyshev cosine approximation using the BSGS (Baby-Step Giant-Step) algorithm [39] and the Paterson-Stockmeyer method, followed by several iterations of the double angle formula to simulate a larger approximation range at the cost of increased level usage [37].

For the CoeffToSlot and SlotToCoeffs steps, FIDESlib generalizes them into a single routine. These steps perform homomorphic encoding and decoding through the homomorphic application of a Discrete Fourier Transform (DFT). To improve efficiency at the cost of increasing multiplicative levels, the DFT plaintext matrix is split into several block matrices with higher sparsity [44]. Each ciphertext-vector times plaintext-matrix multiplication is then performed using a BSGS algorithm [42], which reduces the number of required rotations and leverages the hoisted rotation optimization. Unlike OpenFHE, FIDESlib does not perform ModDown hoisting, this reduces the size of precomputed plaintexts and minimizes device memory usage.

## IV. EVALUATION

### A. Methodology

We have used Google Test to perform unit tests, and also integration tests with OpenFHE, developing a test suite that has been executed successfully using the following NVIDIA GPUs: 4060Ti, A4500, V100, and 4090 (details in Table IV). We compare the efficiency of FIDESlib with OpenFHE (baseline), an AVX-512 optimized OpenFHE implementation lever-

aging Intel's HEXL [12], and the current leading open-source GPU implementation, Phantom.

In order to comprehensively evaluate all possible scenarios, all tests are executed for different CKKS parameter sets ($[N, L, \Delta, dnum]$). Tests for specific operations such as CiphertextSquaring are further specified with additional parameters, such as limb-batching configuration or ciphertext level. For the performance evaluation of FIDESlib, we opted for Google Benchmark, a benchmarking tool that offers a similar interface to Google Test and allows for easy customization to collect a wide range of performance metrics. It also allows for the parametrization of benchmarks, enabling us to easily obtain performance information for FIDESlib using different parameter sets on various GPUs.

### B. Results

We evaluate the performance of every public API operation in FIDESlib, as well as some internal operations, to obtain comprehensive data on all implemented functionality. Unless otherwise specified, $[N, L, \Delta, dnum] := [2^{16}, 29, 59, 4]$.

The main performance results are presented in Table V. We show FIDESlib's operations execution time for a single maximum level ciphertext ($\ell$=29). We compare against a baseline single-threaded OpenFHE installation, multi-threaded and HEXL enabled OpenFHE, and the Phantom FHE open-source GPU library. The CPU tests are run on an AMD Ryzen 9 7900 (12-core, SMT enabled, AVX-512 enabled) with DDR5 memory at 5200 MT/s. All GPU tests are run on the same Nvidia RTX 4090 system. As we can see, FIDESlib achieves the best performance on every operation. Notably, HMult is more than $100\times$ faster than a multi-threaded CPU implementation, while Rescale is more than $30\times$ faster.

A performant and balanced (i)NTT implementation is crucial to achieve fast and reliable FHE performance on a variety of platforms. Figure 4 shows a comparison of the time per limb on a range of limb working sets, in terms of size. For context, a ModUp operation involves computing up to $(dnum + 1) * L$ (around 150 for bootstrappable parameters) (i)NTT operations. Our implementation shows greater memory bandwidth efficiency as the working set grows.

Since plaintext multiplication followed by a rescaling operation is a common pattern in CKKS, we evaluate this sequence in Figure 5. As observed, the time complexity is nearly linear with the number of RNS limbs involved. Some slight non-linearity is observed with the RTX 4060 Ti GPU, as the working set size starts to fit into its 32-MB L2 cache below the 20-limb data point.

Figure 6 illustrates the performance of the HMult operation across various architectures. The limb batch parameter employed on each platform is the one that yielded the highest performance. We showcase how performance varies depending on the amount of used levels of the ciphertext for a given parameter set. The Hybrid-Key-Switching technique causes a noticeable speedup each time a "digit" is dropped (as ciphertext multiplicative levels are used).

| Compute Platform | Frequency | CPU Cores or SMs | 32b INT TOPS | Private Data Cache | Shared Cache | DRAM Size | Bandwidth |
|---|---|---|---|---|---|---|---|
| CPU: Ryzen 9 7900 | 3.70 GHz | 12 | 2.13 | 1056 KB | 64 MB | 64 GB | 81 GB/s |
| GPU: RTX 4060 Ti | 2.31 GHz | 34 | 11.03 | 128 KB | 32 MB | 16 GB | 288 GB/s |
| GPU: RTX A4500 | 1.05 GHz | 56 | 11.83 | 128 KB | 6 MB | 20 GB | 640 GB/s |
| GPU: V100 | 1.25 GHz | 80 | 14.13 | 128 KB | 6 MB | 16 GB | 897 GB/s |
| GPU: RTX 4090 | 2.24 GHz | 128 | 41.29 | 128 KB | 72 MB | 24 GB | 1 TB/s |

TABLE IV: Specifications of the CPU and GPU compute platforms used in our experiments.

| | OpenFHE (Baseline) | OpenFHE (Intel HEXL. 24 threads) | | Phantom (RTX 4090) | | **FIDESlib** (RTX 4090) | |
|---|---|---|---|---|---|---|---|
| **ScalarAdd** | 1.28 ms | 106.00 $\mu$s | 12.06$\times$ | N/A | | **16.63** $\mu$s | 76.89$\times$ |
| **PtAdd** | 5.26 ms | 5.80 ms | 0.90$\times$ | 20.64 $\mu$s | 254.87$\times$ | **17.79** $\mu$s | 295.65$\times$ |
| **HAdd** | 7.84 ms | 8.39 ms | 0.93$\times$ | 82.66 $\mu$s | 94.86$\times$ | **50.70** $\mu$s | 154.67$\times$ |
| **ScalarMult** | 4.34 ms | 225.00 $\mu$s | 19.29$\times$ | N/A | | **44.15** $\mu$s | 98.32$\times$ |
| **PtMult** | 10.14 ms | 5.32 ms | 1.90$\times$ | 31.91 $\mu$s | 317.25$\times$ | **21.74** $\mu$s | 465.72$\times$ |
| **Rescale** | 50.80 ms | 4.92 ms | 10.31$\times$ | 224.58 $\mu$s | 226.21$\times$ | **156.11** $\mu$s | 325.44$\times$ |
| **HRotate\*** | 370.71 ms | 105.30 ms | 3.52$\times$ | 1.139 ms | 325.60$\times$ | **1.107** ms | 334.90$\times$ |
| **HMult** | 406.24 ms | 151.58 ms | 2.62$\times$ | 1.220 ms | 332.99$\times$ | **1.084** ms | 374.61$\times$ |

TABLE V: Performance comparison of CKKS primitives. *Conjugate is implemented exactly the same as HRotate.



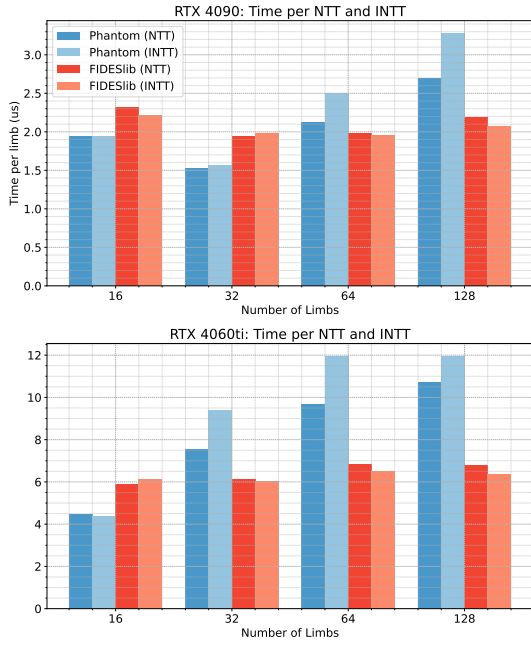Fig. 4: (i)NTT performance comparison (lower is better).



Fig. 5: $PtMult$+$Rescale$. Best limb batch.



Fig. 6: $HMult$. Best limb batch.

Figure 7 highlights the impact of the limb batch parameter on the performance of the HMult operation. As GPU throughput increases, the best performing configuration is found at a higher limb batch.

Performance of the HMult operation is shown for different parameter sets in Figure 8, showcasing FIDESlib's versatility and robustness. On smaller parameter sets, the workload is small and it is bottlenecked by kernel latency, so higher frequency platforms (RTX 4060 ti, RTX 4090) see improved performance over throughput oriented designs (V100). Furthermore, as key-switching key sizes vary (2.3MB, 7.7MB, 20MB, 152MB, and 360MB, respectively), some devices ben-
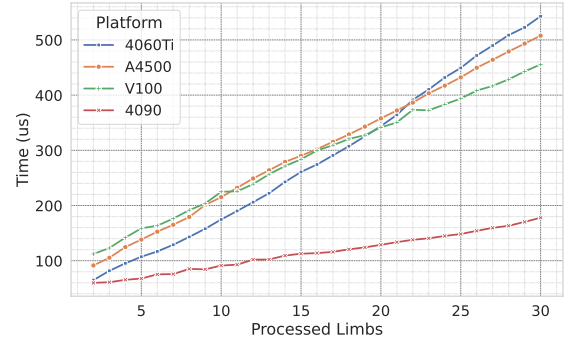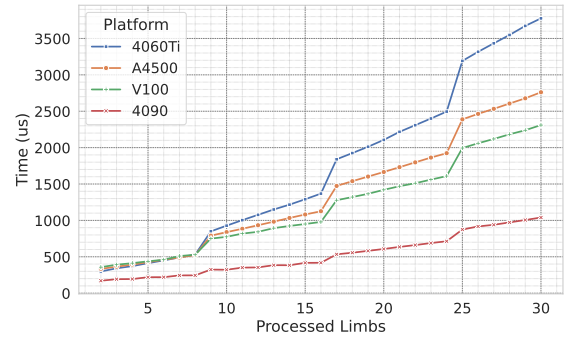
efit from higher L2 cache hit-rate at some parameter levels, improving performance.

Bootstrapping performance and throughput, compared to our baseline OpenFHE implementation and the Intel-optimized OpenFHE HEXL version, is presented in Table VI. The execution times are comparable as the achieved message precision is equal across all implementations. Execution time increases with the number of slots (encoded values), though is amortized over time as we increase the number of slots. Performance is
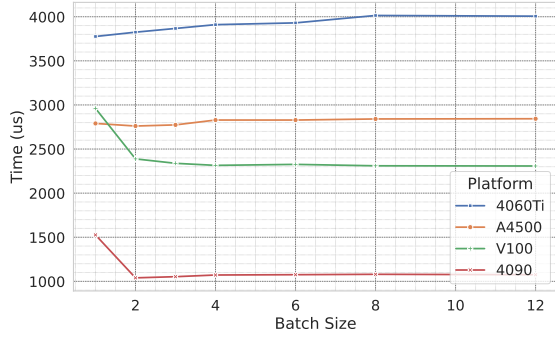
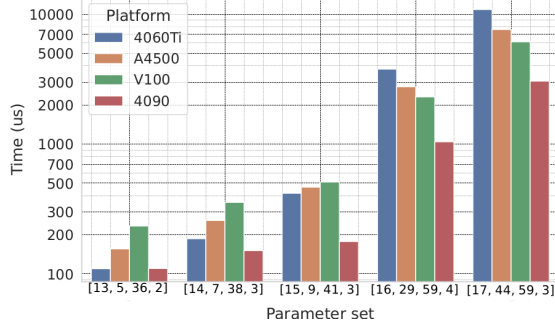Fig. 7: $HMult$ with maximum levels. Limb batch comparison.



Fig. 8: $HMult$ with maximum levels. Parameter set $[log(N), L, \Delta, dnum]$ comparison.

| Slots | Levels | OpenFHE | HEXL-24 threads | **FIDESlib**-RTX 4090 |
|---|---|---|---|---|
| 64 | 13 | 18 224 $T$ | 5 204 (3.50×) | **73.5** (248×) |
| | | 21 904 $A$ | 6 255 | **88.3** |
| 512 | 11 | 18 268 $T$ | 7 781 (2.35×) | **93.3** (196×) |
| | | 3 243 $A$ | 1 381 | **16.6** |
| 16384 | 9 | 20 079 $T$ | 9 281 (2.16×) | **112** (179×) |
| | | 136 $A$ | 62.9 | **0.761** |
| 32768 | 9 | 28 635 $T$ | 12 185 (2.35×) | **146** (196×) |
| | | 97.1 $A$ | 41.3 | **0.496** |

TABLE VI: Bootstrapping performance. $Slots$=Number of encoded values; $Levels$=Levels remaining after bootstrapping; $T$=Time in ms; $A$=Amortized time in $\frac{\mu s}{\text{slots} \cdot \text{levels}}$.

| Time (ms) | OpenFHE | HEXL - 24 threads | FIDESlib - RTX 4090 |
|---|---|---|---|
| Iteration | 1,555 | 448 (3.47×) | **23** (67.61×) |
| Iteration + Bootstrap | 16,233 | 7,233 (2.24×) | **169** (96.05×) |

TABLE VII: LR performance. Time in ms (and speedup).

over 70× faster than HEXL-enabled OpenFHE.

To evaluate our library's robustness and performance under realistic workloads, we implemented a modified Logistic Regression (LR) algorithm based on the approach in [51]. We trained the model on a dataset of 45,000 loan eligibility samples using mini-batch standard gradient descent with ciphertexts containing 1,024 samples. Although each data sample had 25 parameters after encoding, we opted to align our data to the next power of two boundary, 32, to optimize the rotations needed on each iteration. Table VII shows the performance comparison of the same training algorithm both on OpenFHE and FIDESlib. We bootstrap every LR iteration while employing parameters $[N, L, \Delta, dnum] := [2^{16}, 26, 59, 4]$.

| Features | [17] | [18] | [19] | [20] | [15] | [16] | [23] | [22] | FIDESlib |
|---|---|---|---|---|---|---|---|---|---|
| Open Source | | ✓ | | ✓ | ✓ | | ✓ | | ✓ |
| Published | | | ✓ | | ✓ | ✓ | | ✓ | ✓ |
| Bootstrapping | ✓ | | ✓ | | | ✓ | | ✓ | ✓ |
| OpenFHE Inter. | | | | | | | | | ✓ |
| Benchmarks | | | ✓ | | | ✓ | | ✓ | $LR$ |
| Microbench. | | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | ✓ |
| Unit Tests | ✓ | | ✓ | | | | | | ✓ |
| Integration Tests | | | | | | | | | ✓ |
| Multi-GPU | | | | ✓ | | | ✓ | | $WIP$ |

TABLE VIII: Qualitative comparison of GPU CKKS libraries. $LR$=Logistic Regression; $WIP$=Work in progress.

## V. RELATED WORK

Table VIII shows a qualitative comparison of the state-of-the-art GPU-based CKKS libraries, while consider a range of features. Phantom [15] is the most efficient open-source CKKS library for a GPU backend up to date, and is why we have used it as our GPU baseline for performance analysis in the previous section. However, it lacks several CKKS operations such as ScalarAdd, ScalarMult, HSquare, and bootstrapping. Libraries such as HEaaN [17], Over100x [19], Cheddar [16], and TensorFHE [22] offer feature-complete CKKS implementations, including bootstrapping. Nonetheless, the lack of publicly available source code hinders reproducibility and further research. Older versions of HEaaN and Over100x are publicly available; however, to the best of our knowledge, the open-source HEaaN code only supports a CPU backend, and the Over100x GPU backend is incomplete, preventing full reproduction of the reported experiments and limiting support for other CKKS parameter settings.

In contrast, FIDESlib is the first fully-fledged open-source CKKS library designed for GPUs, delivering better performance than Phantom. FIDESlib is designed to be fully compatible with OpenFHE [34], the widely trusted industry standard for FHE. FIDESlib is the only library that ensures complete interoperability with OpenFHE, accelerating server-side CKKS operations, while offloading encoding and encryption tasks to OpenFHE-based clients.

Robust testing of library functionality is lacking in this field; only HEonGPU [18], Troy-Nova [20] (both are available for download, but have not been published in peer-reviewed venues), and OpenFHE [34] offer some unit tests. In contrast, FIDESlib features an extensive suite of unit tests that, combined with our extensive set of microbenchmarks, and our Logistic Regression benchmark, thoroughly validate all implemented functionality. Additionally, our new library represents, to our knowledge, the first instance of validation enhancement using integration tests that compare FIDESlib's output with OpenFHE results. Moreover, leveraging Google Test and Google Benchmark makes implementing continuous integration in our library straightforward.

Although Liberate-FHE [23] and Troy-Nova (both unpublished, non-peer-reviewed libraries) support multi-GPU backends for CKKS, Troy-Nova imposes limitations, such as requiring manual transfer of objects between devices. Liberate-FHE does support multi-GPU by dividing the workloads at the RNS level, but their implementation is an order of magnitude slower than other state-of-the-art implementations. FIDESlib's

architecture will allow for the implementation of a more scalable multi-GPU solution, built upon a more performant baseline.

## VI. Conclusions and Future Work

This paper introduces FIDESlib, the first open-source CKKS library with full functionality optimized for a GPU backend, including bootstrapping, and featuring robust security from OpenFHE's client-side operations. Compared to Phantom, the current best-performing alternative, FIDESlib delivers better scalability and performance across all shared operations on four tested GPU platforms. For the most time-consuming and challenging bootstrapping operation, which is not supported in Phantom, FIDESlib achieves up to $227.8\times$ and $74.4\times$ speedups over OpenFHE and AVX-optimized OpenFHE, respectively. FIDESlib is built with modern software practices, supporting extensibility and broad adoption. Our roadmap includes AMD GPU support via HIP [52], multi-GPU backends for NVIDIA and AMD, and porting more real-world MLaaS cloud-based workloads.

## References

[1] M. Ribeiro, K. Grolinger, and M. A. Capretz, "Mlaas: Machine learning as a service," in *2015 IEEE 14th international conference on machine learning and applications (ICMLA)*. IEEE, 2015, pp. 896–902.

[2] A. Mishra, *Machine learning in the AWS cloud: Add intelligence to applications with Amazon Sagemaker and Amazon Rekognition*. John Wiley & Sons, 2019.

[3] S. R. Sukhdeve and S. S. Sukhdeve, *Google Cloud Platform for Data Science*. Apress Berkeley, CA, 2023.

[4] C. Körner and K. Waaijer, *Mastering Azure Machine Learning: Perform large-scale end-to-end advanced machine learning in the cloud with Microsoft Azure Machine Learning*. Packt Publishing Ltd, 2020.

[5] IBM, "Cost of a Data Breach Report," 2024, https://www.ibm.com/reports/data-breach.

[6] S. Lloyd-Jones and K. Manwaring, "First steps to quantum resilience: Identifying 'broken concepts' in australia's national security laws," *ANU Journal of Law and Technology*, September 2024, forthcoming. [Online]. Available: https://ssrn.com/abstract=4976322

[7] National Institute of Standards and Technology (NIST), "FIPS 203: Module-Lattice-Based Key-Encapsulation Mechanism Standard," August 2024, https://csrc.nist.gov/pubs/fips/203/final.

[8] National Institute of Standards and Technology, "FIPS 204: Module-Lattice-Based Digital Signature Standard," August 2024, https://csrc.nist.gov/pubs/fips/204/final.

[9] F. R. T. D. J. of the United States Government, "Announcing Issuance of Federal Information Processing Standards," August 2024, https://www.federalregister.gov/d/2024-17956.

[10] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I 23*. Springer, 2017, pp. 409–437.

[11] W. Jung, E. Lee, S. Kim, J. Kim, N. Kim, K. Lee, C. Min, J. H. Cheon, and J. H. Ahn, "Accelerating fully homomorphic encryption through architecture-centric analysis and optimization," *IEEE Access*, vol. 9, pp. 98 772–98 789, 2021.

[12] F. Boemer, S. Kim, G. Seifu, F. D. M. de Souza, and V. Gopal, "Intel hexl: Accelerating homomorphic encryption with intel avx512-ifma52," 2021. [Online]. Available: https://arxiv.org/abs/2103.16400

[13] Y. Gong, X. Chang, J. Mišić, V. B. Mišić, J. Wang, and H. Zhu, "Practical solutions in fully homomorphic encryption: a survey analyzing existing acceleration methods," *Cybersecurity*, vol. 7, no. 1, 3 2024. [Online]. Available: https://doi.org/10.1186/s42400-023-00187-4

[14] Intel, "Why Data Center GPUs Are Essential to Innovation," https://www.intel.com/content/www/us/en/products/docs/discrete-gpus/data-center-gpu/what-is-data-center-gpu.html, 2024, [Online; accessed 13-December-2024].

[15] H. Yang, S. Shen, W. Dai, L. Zhou, Z. Liu, and Y. Zhao, "Phantom: a cuda-accelerated word-wise homomorphic encryption library," *IEEE Transactions on Dependable and Secure Computing*, 2024.

[16] J. Kim, W. Choi, and J. H. Ahn, "Cheddar: A swift fully homomorphic encryption library for cuda gpus," *arXiv preprint arXiv:2407.13055*, 2024.

[17] C. INC., "Heaan private ai homomorphic encryption library. fully homomorphic encryption with ckks scheme, faster bootstrapping and gpu-acceleration." 2022, https://heaan.it/.

[18] A. Şah Özcan and E. Savaş, "HEonGPU: a GPU-based fully homomorphic encryption library 1.0," Cryptology ePrint Archive, Paper 2024/1543, 2024. [Online]. Available: https://eprint.iacr.org/2024/1543

[19] W. Jung, S. Kim, J. H. Ahn, J. H. Cheon, and Y. Lee, "Over 100x faster bootstrapping in fully homomorphic encryption through memory-centric optimization with gpus," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 114–148, 2021.

[20] lightbulb128, "The new implementation of troy, a cuda based gpu parallelized implementation of rlwe homomorphic encryption schemes." 2023, https://github.com/lightbulb128/troy-nova.

[21] Z. Wang, P. Li, R. Hou, Z. Li, J. Cao, X. Wang, and D. Meng, "Hebooster: An efficient polynomial arithmetic acceleration on gpus for fully homomorphic encryption," *IEEE Transactions on Parallel and Distributed Systems*, vol. 34, no. 4, pp. 1067–1081, 2023.

[22] S. Fan, Z. Wang, W. Xu, R. Hou, D. Meng, and M. Zhang, "Tensorfhe: Achieving practical computation on encrypted data using gpgpu," in *2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2023, pp. 922–934.

[23] DESILO, "Liberate.FHE: A New FHE Library for Bridging the Gap between Theory and Practice with a Focus on Performance and Accuracy," 2023, https://github.com/Desilo/liberate-fhe, https://github.com/FHE-org/fhe-org.github.io/files/14896471/07-Kim.pdf.

[24] Google, "Google test," 2024, accessed: 2024-12-17. [Online]. Available: https://github.com/google/googletest

[25] ——, "Google benchmark," 2024, accessed: 2024-12-17. [Online]. Available: https://github.com/google/benchmark

[26] C. Gentry, "Fully homomorphic encryption using ideal lattices," in *Proceedings of the forty-first annual ACM symposium on Theory of computing*, 2009, pp. 169–178.

[27] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(leveled) fully homomorphic encryption without bootstrapping," *ACM Transactions on Computation Theory (TOCT)*, vol. 6, no. 3, pp. 1–36, 2014.

[28] J. Fan and F. Vercauteren, "Somewhat practical fully homomorphic encryption," *Cryptology ePrint Archive*, 2012.

[29] I. Chillotti, N. Gama, M. Georgieva, and M. Izabachène, "Tfhe: fast fully homomorphic encryption over the torus," *Journal of Cryptology*, vol. 33, no. 1, pp. 34–91, 2020.

[30] L. Ducas and D. Micciancio, "Fhew: bootstrapping homomorphic encryption in less than a second," in *Annual international conference on the theory and applications of cryptographic techniques*. Springer, 2015, pp. 617–640.

[31] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and*

*Applications of Cryptology and Information Security, Hong Kong, China, December 3-7, 2017, Proceedings, Part I 23*. Springer, 2017, pp. 409–437.

[32] S. Halevi and V. Shoup, "Design and implementation of helib: a homomorphic encryption library," *Cryptology ePrint Archive*, 2020.

[33] K. Laine, "Simple encrypted arithmetic library 2.3. 1," *Microsoft Research https://www. microsoft. com/en-us/research/uploads/prod/2017/11/sealmanual-2-3-1. pdf*, 2017.

[34] A. Al Badawi, J. Bates, F. Bergamaschi, D. B. Cousins, S. Erabelli, N. Genise, S. Halevi, H. Hunt, A. Kim, Y. Lee *et al.*, "Openfhe: Open-source fully homomorphic encryption library," in *proceedings of the 10th workshop on encrypted computing & applied homomorphic cryptography*, 2022, pp. 53–63.

[35] J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song, "A full rns variant of approximate homomorphic encryption," in *Selected Areas in Cryptography–SAC 2018: 25th International Conference, Calgary, AB, Canada, August 15–17, 2018, Revised Selected Papers 25*. Springer, 2019, pp. 347–368.

[36] A. Kim, A. Papadimitriou, and Y. Polyakov, "Approximate homomorphic encryption with reduced approximation error," in *Cryptographers' Track at the RSA Conference*. Springer, 2022, pp. 120–144.

[37] K. Han and D. Ki, "Better bootstrapping for approximate homomorphic encryption," in *Cryptographers' Track at the RSA Conference*. Springer, 2020, pp. 364–390.

[38] J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song, "Bootstrapping for approximate homomorphic encryption," in *Advances in Cryptology– EUROCRYPT 2018: 37th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Tel Aviv, Israel, April 29-May 3, 2018 Proceedings, Part I 37*. Springer, 2018, pp. 360–384.

[39] S. Halevi and V. Shoup, "Faster homomorphic linear transformations in helib," in *Annual International Cryptology Conference*. Springer, 2018, pp. 93–120.

[40] J. H. Cheon, K. Han, and M. Hhan, "Faster homomorphic discrete fourier transforms and improved fhe bootstrapping," *Cryptology ePrint Archive*, 2018.

[41] J.-W. Lee, E. Lee, Y. Lee, Y.-S. Kim, and J.-S. No, "High-precision bootstrapping of rns-ckks homomorphic encryption using optimal minimax polynomial approximation and inverse sine function," in *Advances in Cryptology–EUROCRYPT 2021: 40th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Zagreb, Croatia, October 17–21, 2021, Proceedings, Part I 40*. Springer, 2021, pp. 618–647.

[42] J.-P. Bossuat, C. Mouchet, J. Troncoso-Pastoriza, and J.-P. Hubaux, "Efficient bootstrapping for approximate homomorphic encryption with non-sparse keys," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2021, pp. 587–617.

[43] J.-P. Bossuat, J. Troncoso-Pastoriza, and J.-P. Hubaux, "Bootstrapping for approximate homomorphic encryption with negligible failure-probability by using sparse-secret encapsulation," in *International Conference on Applied Cryptography and Network Security*. Springer, 2022, pp. 521–541.

[44] H. Chen, I. Chillotti, and Y. Song, "Improved bootstrapping for approximate homomorphic encryption," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2019, pp. 34–54.

[45] R. Agrawal, L. De Castro, C. Juvekar, A. Chandrakasan, V. Vaikuntanathan, and A. Joshi, "Mad: Memory-aware design techniques for accelerating fully homomorphic encryption," in *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, ser. MICRO '23, New York, NY, USA, 2023. [Online]. Available: https://doi.org/10.1145/3613424.3614302

[46] J. Kim, G. Lee, S. Kim, G. Sohn, M. Rhu, J. Kim, and J. H. Ahn, "Ark: Fully homomorphic encryption accelerator with runtime data generation and inter-operation key reuse," in *2022 55th IEEE/ACM International Symposium on Microarchitecture (MICRO)*. IEEE, Oct. 2022, p. 1237–1254. [Online]. Available: http://dx.doi.org/10.1109/MICRO56248.2022.00086

[47] L. de Castro, R. Agrawal, R. Yazicigil, A. Chandrakasan, V. Vaikuntanathan, C. Juvekar, and A. Joshi, "Does fully homomorphic encryption need compute acceleration?" 2021. [Online]. Available: https://arxiv.org/abs/2112.06396

[48] R. Agrawal and A. Joshi, *On Architecting Fully Homomorphic Encryption-based Computing Systems*, 01 2023.

[49] K. Shivdikar, Y. Bao, R. Agrawal, M. Shen, G. Jonatan, E. Mora, A. Ingare, N. Livesay, J. L. Abellán, J. Kim *et al.*, "Gme: Gpu-based microarchitectural extensions to accelerate homomorphic encryption," in *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, 2023, pp. 670–684.

[50] K. Shivdikar, G. Jonatan, E. Mora, N. Livesay, R. Agrawal, A. Joshi, J. L. Abellán, J. Kim, and D. Kaeli, "Accelerating polynomial multiplication for homomorphic encryption on gpus," in *2022 IEEE International Symposium on Secure and Private Execution Environment Design (SEED)*. IEEE, 2022, pp. 61–72.

[51] K. Han, S. Hong, J. H. Cheon, and D. Park, "Logistic regression on homomorphic encrypted data at scale," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 9466–9471, Jul. 2019. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/5000

[52] Y. Sun, T. Baruah, and D. Kaeli, "Accelerated Computing with HIP," https://community.amd.com/t5/instinct-accelerators/accelerated-computing-with-hip-textbook/ba-p/620840, 2022, ingram Spark Publishing.