

# HMC-FHE: A Heterogeneous Near Data Processing Framework for Homomorphic Encryption

Zehao Chen<sup>1b</sup>, Zhining Cao, Zhaoyan Shen<sup>1b</sup>, and Lei Ju<sup>1b</sup>

**Abstract**—Fully homomorphic encryption (FHE) offers a promising solution to ensure data privacy by enabling computations directly on encrypted data. However, its notorious performance degradation severely limits the practical application, due to the explosion of both the ciphertext volume and computation. In this article, leveraging the diversity of computing power and memory bandwidth requirements of FHE operations, we present HMC-FHE, a robust acceleration framework that combines both GPU and hybrid memory cube (HMC) processing engines to accelerate FHE applications cooperatively. HMC-FHE incorporates four key hardware/software co-design techniques: 1) a fine-grained kernel offloading mechanism to efficiently offload FHE operations to relevant processing engines; 2) a ciphertext partitioning scheme to minimize data transfer across decentralized HMC processing engines; 3) an FHE operation pipeline scheme to facilitate pipelined execution between GPU and HMC engines; and 4) a kernel tuning scheme to guarantee the parallelism of GPU and HMC engines. We demonstrate that the GPU-HMC architecture with proper resource management serves as a promising acceleration scheme for memory-intensive FHE operations. Compared with the state-of-the-art GPU-based acceleration scheme, the proposed framework achieves up to 2.65× performance gains and reduces 1.81× energy consumption with the same peak computation capacity.

**Index Terms**—Accelerator, homomorphic encryption, near memory processing.

## I. INTRODUCTION

FULLY homomorphic encryption (FHE) is an encryption scheme that enables computations to be performed directly on encrypted data, where no decryption for

intermediate steps during the computation is required. With the ability to perform arithmetic operations without revealing user data privacy, FHE becomes one of the most promising privacy protection techniques and has been gradually deployed in application scenarios, including encrypted databases [1] and machine learning [2].

Nonetheless, the concerning computational speed of FHE remains a major obstacle hindering its rapid advancement. Since the method of constructing the ciphertext space results in a substantial expansion of computation and data scale, the computing speed of ciphertext experiences a drastic degradation of several orders of magnitudes ( $10^3 \sim 10^6$ ) compared to plaintext [3]. To solve this issue, numerous studies [4], [5], [6], [7], [8], [9], [10], [11] have focused on building domain-specific hardware accelerators to enhance the performance of FHE operations through resource reuse and increased parallelism. Meanwhile, researchers have also explored deploying FHE operations on GPU platforms to support privacy protection applications in various scenarios comprehensively [12], [13], [14], [15]. These efforts have yielded a substantial enhancement in the speed of ciphertext computation, reducing the performance gap by 2~3 orders of magnitude.

In the design of FHE accelerators, due to the memory-intensive nature, the demand for memory capacity and bandwidth far surpasses the computational requisites. For instance, ASIC schemes typically use large on-chip storage [7], [16]. Meanwhile, FPGA (or GPU)-based accelerator designs identify on-chip BRAM (or Shared Memory) capacity as the performance bottleneck [4], [12]. In this work, we further conduct extensive quantitative analysis on the arithmetic intensity (AI) of each FHE operation on GPU devices. Our analysis results further reinforce the significance of prioritizing the alleviation of bandwidth constraints as the foremost endeavor in accelerating FHE operations.

The memory-intensive characteristics of FHE naturally make hybrid memory cube (HMC) a natural choice for FHE acceleration. In particular, the GPU-HMC architecture (as demonstrated in previous work [17], [18], [19], [20], [21]) as shown in Fig. 1 is a promising candidate which is composed of a centralized GPU device and multiple auxiliary HMC processing engines. By encapsulating various computing logic for computing while having the characteristics of high-internal bandwidth, low latency, and low-power consumption, the GPU-HMC architecture offers a balanced tradeoff between computing power and bandwidth resources for individual FHE operations. Nonetheless, relying solely on this architectural

Manuscript received 10 August 2024; accepted 10 August 2024. This work was supported in part by the Natural Science Foundation of China under Grant 62372272; in part by the Department of Science and Technology of Shandong Province under Grant SYS202201; in part by the Quan Cheng Laboratory under Grant QCLZD202302 and Grant CCF-AFSG RF20220011; in part by the Taishan Scholars Program under Grant tsqn202211281; in part by the Key Research and Development Program of Shandong Province under Grant 2023CXPT002; and in part by the Open Project Program of Wuhan National Laboratory for Optoelectronics under Grant 2022WNLOKF018. This article was presented at the International Conference on Hardware/Software Codesign and System Synthesis (CODES + ISSS) 2024 and appeared as part of the ESWEEK-TCAD Special Issue. This article was recommended by Associate Editor S. Dailey. (Corresponding author: Lei Ju.)

Zehao Chen is with the School of Cyber Science and Technology, Shandong University, Qingdao 266237, China, and also with Quan Cheng Laboratory, Jinan 250012, China.

Zhining Cao is with the School of Computer Science and Technology, Shandong University, Qingdao 266237, China.

Zhaoyan Shen is with the School of Computer Science and Technology, Shandong University, Qingdao 266237, China, and also with the Wuhan National Laboratory for Optoelectronics, Huazhong University of Science and Technology, Wuhan 430074, China.

Lei Ju is with Quan Cheng Laboratory, Jinan 250012, China (e-mail: julei@sdu.edu.cn).

Digital Object Identifier 10.1109/TCAD.2024.3447212

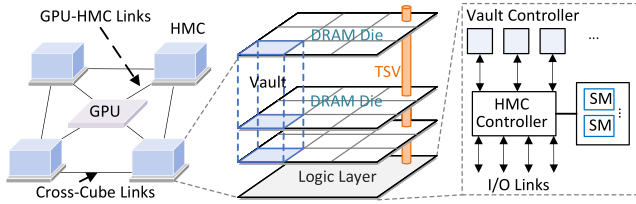


Fig. 1. Typical GPU-HMC architecture [17], [18].

78 setup leads to suboptimal overall performance enhancement,  
79 and several challenges must be effectively addressed with an  
80 automatic design flow.

81 *Challenge 1:* Different FHE operations exhibit varying com-  
82 putation and memory intensity. Thus, exquisitely offloading  
83 FHE operations to their preferred processing engines (GPU or  
84 HMCs) becomes a primary problem to address.

85 *Challenge 2:* To execute an FHE operation on distributed  
86 HMC processing engines in a parallel fashion, it is essential to  
87 reduce the inter-HMC data transfer between parallel executed  
88 subtasks.

89 *Challenge 3:* It is critical to utilize the parallel compu-  
90 tation capacity between GPU and HMCs for overall system  
91 performance, which necessitates sophisticated scheduling to  
92 achieve inter- and intra-operation parallelism.

93 In this article, to overcome the above challenges, we  
94 propose a design flow framework for GPU-HMC-based FHE  
95 acceleration, which consists of a series of hardware/software  
96 co-designs. First, we introduce an offloading mechanism  
97 for FHE operations, which categorizes the operations into  
98 reusable basic kernels and determines the affinity for the  
99 GPU or HMC engines based on the AI. Moreover, we  
100 propose a ciphertext partitioning scheme by decoupling the  
101 structure of ciphertext polynomials, which ensures efficient  
102 global memory bandwidth utilization while achieving paral-  
103 lelism and load balancing. Then, we undertake a thorough  
104 analysis of data dependency and propose an interoperation  
105 pipeline scheme to facilitate parallel execution between GPU  
106 and HMC engines. Finally, for HE operations outside the  
107 above-mentioned pipeline stages, we propose a fine-grained  
108 intraoperation tuning scheme to further balance the workload  
109 between GPU and HMC engines.

110 The contributions of this work are summarized as follows.

- 111 1) To the best of the authors knowledge, this is a pioneer  
112 work that introduces FHE operation acceleration with  
113 heterogeneous GPU-HMC architecture. It unveils that  
114 GPU-HMC architecture is a promising design choice  
115 for FHE acceleration given the distinct computation and  
116 memory intensity between FHE operations.
- 117 2) This article proposes a design flow framework which  
118 automatically performs FHE operation scheduling and  
119 data allocation on the GPU-HMC architecture. The  
120 framework synergistically integrates a set of hard-  
121 ware/software co-design techniques to achieve working  
122 balancing and high-memory bandwidth utilization.
- 123 3) The experimental results on practical FHE workloads  
124 show up-to  $2.65\times$  speedup and up-to  $1.81\times$  energy  
125 efficiency (EE), compared to the state-of-the-art GPU-  
126 based accelerator design. Meanwhile, the scalability of

the proposed GPU-HMC architecture has also been  
evaluated.

The remainder of this article is organized as fol-  
lows. Section II introduces the background and motivation.  
Section III gives an overview of HMC-FHE framework.  
Section IV describe the detail techniques of HMC-FHE.  
Section V evaluates HMC-FHE. Sections VI and VII discuss  
the related work and conclusion.

## II. BACKGROUND AND MOTIVATION

In this section, we first introduce the basic of FHE with a  
typical algorithm CKKS [22]. Next, we briefly introduce the  
structure of hybrid memory cube (HMC) processing engines.  
Lastly, we provide the motivation of this work.

### A. FHE Schemes

FHE, exemplified by schemes, such as CKKS [22], BFV [23],  
and BGV [24], enables arbitrary computations on ciphertexts  
without the need for decryption. This article primarily con-  
centrates on the CKKS scheme due to its unique capability  
to handle plaintext inputs comprising arbitrary fixed-point real  
numbers, which has facilitated its application across a broader  
range of fields. Note that the proposed techniques and schemes  
are equally applicable to other FHE schemes, e.g., BFV, and  
BGV.

*Ciphertext Structure in CKKS:* In the CKKS scheme, each  
batch of  $N/2$  real numbers is encoded and encrypted as  
a pair of polynomials with degree  $N$  in the ring  $R_Q =$   
 $Z_Q[X]/(X^N + 1)$ , where  $Q$  is a prime number with hundreds or  
thousands of bits that relate directly to the ciphertext spaces.  
Each polynomial consists of  $N$  coefficients, which are integers  
modulo by  $Q$ . To enable efficient modulo computations of  
the wide  $Q$  and coefficients, the residue number system  
(RNS) is presented to convert these wide coefficients into  
 $L$  residue polynomials with machine-word-length coefficients.  
Consequently, the ciphertext can be initially represented as a  
pair of 2-D matrices with width  $L$  and depth  $N$ .

*CKKS Basic Kernel:* The CKKS algorithm is built upon  
a set of basic kernels. All basic kernels are summarized as  
follows.

- 1) *Addition/Subtraction (Add/Sub)*, which execute element-  
wise operations to add or subtract two polynomials.
- 2) *Tensor Product (TensorP)*, which conducts element-wise  
product of two polynomials.
- 3) *Number Theory Transformation (NTT and iNTT)*, which  
enables the conversion between polynomials represented  
by coefficients and those represented by point-values,  
significantly speeding up polynomial multiplication.
- 4) *Fast Basic Conversion (Conv)*, which converts the RNS  
basis based on  $B = \{p_0, p_1, \dots, p_{K-1}\}$  and  $C =$   
 $\{q_0, q_1, \dots, q_l\}$ ,  $0 \leq l \leq L$ . We refer to the operations  
occurring during modulus increase or modulus decrease as  
*ConvUp (ConvU)* and *ConvDown (ConvD)*, respectively.
- 5) *Inner Product (InnerP)*, which integrates several  
*TensorP* kernels along with a single *Add* kernel to  
accumulate the results of *TensorP*.

TABLE I  
DIFFERENT CKKS OPERATIONS

APIs	Description	Basic Kernels
CCAdd	Add two ciphertexts.	Add
PCMult	Multiply a plaintext with a ciphertext.	TensorP
CCMult	Multiply a ciphertext with a ciphertext.	TensorP, NTT, iNTT, Conv, InnerP
Rotate	Rotate the ciphertexts according to an index.	NTT, iNTT, Conv, Sub, InnerP, AMorph
HRotate	Batch processing of Rotate operation.	NTT, iNTT, Conv, Sub, InnerP, AMorph

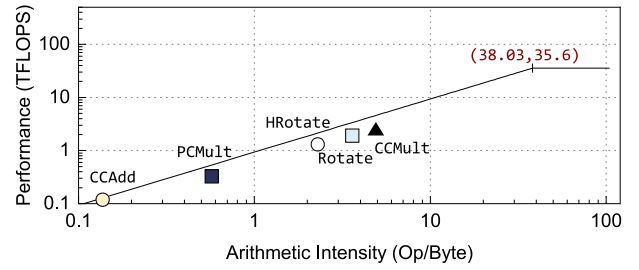


Fig. 2. Roofline model of different CKKS APIs on NVIDIA RTX Geforce 3090.

181 6) *Automorph (AMorph)*, which performs the permutation  
182 operation on the polynomial using an index  $r$ .

183 *CKKS Compound Kernel*: The aforementioned basic kernels  
184 are typically amalgamated into more complex compound  
185 kernels, as outlined below.

186 1) *ModUp and ModDown* adjust the precision of the  
187 ciphertext modulus. *ModUp* enhances precision to  
188 accommodate more sophisticated CKKS operations,  
189 whereas *ModDown* lowers precision to control noise  
190 expansion. Both *ModUp* and *ModDown* encompass  
191 *iNTT*, *ConvU* (for *ModUp*)/*ConvD* (for *ModDown*), and  
192 *NTT* kernels.

193 2) *KeySwitch* facilitates the efficient and secure transition  
194 of ciphertexts encrypted at varying levels. This mech-  
195 anism is crucial in CKKS for operations like *CCMult*,  
196 *Rotate*. Specifically, it consists of *ModUp*, *InnerP*, and  
197 *ModDown* kernels.

198 3) *Rescale* is designed to constrain noise within a specified  
199 range. It is composed of *iNTT*, *NTT*, and *Sub* kernels.

200 4) *Bootstrapping* mitigates the accumulated error through-  
201 out the computational process and resets the ciphertext's  
202 noise level, thereby enabling continued operations on the  
203 ciphertext. This operation encompasses all basic kernels.

204 *CKKS APIs*: By reorganizing these kernels based on spe-  
205 cific rules, CKKS provides a series of APIs for external  
206 applications to implement homomorphic computation. Table I  
207 details the functional description and composition of kernels  
208 for these APIs.

## 209 B. Hybrid Memory Cube

210 Fig. 1 illustrates a heterogeneous GPU-HMC architecture  
211 with four HMCs, and each HMC adopts a 3-D-stacking  
212 architecture within a single package, layering multiple memory  
213 dies atop a logic die through the utilization of through-silicon  
214 via (TSV) technology. Each memory die is subdivided into  
215 several partitions, which are vertically aligned to form vaults,  
216 analogous to conventional memory channels. To enhance par-  
217 allel processing across the vaults in the HMC, a dedicated vault  
218 controller is assigned to each vault. The logic inside an HMC  
219 typically incorporates one or more streaming multiprocessors  
220 (SMs) [25] or bespoke processing units [26], tailored to boost  
221 performance across a variety of applications. Data inside an  
222 HMC can be directly transferred to its logic layer via the  
223 TSV technology, allowing for short-path data transfer and, in  
224 combination with the parallel structure of vaults, achieving

225 high-internal bandwidth and low-access latency, whereas since  
226 the number of processing units in the HMC logic layer is  
227 relatively limited, it offers reduced computational capacity  
228 compared to GPUs, making it generally suitable for memory-  
229 intensive workloads.

230 Moreover, HMC exhibits impressive scalability. Multiple  
231 HMC devices can be interconnected using I/O links on a  
232 silicon interposer, enabling parallel operations, or they can be  
233 connected to GPUs to facilitate joint processing efforts. As  
234 illustrated in Fig. 1, packet-based protocols are utilized for  
235 communications between HMCs or between HMCs and GPUs,  
236 supporting the bidirectional flow of commands and data. It  
237 should be noted that the latency for data transfer between  
238 HMC devices via I/O links is over three times higher than  
239 the latency for local access by the processing units within the  
240 logic layer.

## 241 C. Motivation

242 In this section, we first provide two observations of CKKS  
243 kernels with some preliminary experiments. Then, we discuss  
244 the advantages of a GPU-HMC heterogeneous architecture in  
245 enhancing the performance of FHE operations.

246 1) *Observation 1—CKKS APIs Exhibit a Greater Demand*  
247 *for Memory Bandwidth Compared to Computational*  
248 *Requirements*: The arithmetic logic of the CKKS scheme  
249 is based on the polynomial field of ciphertext, so it  
250 inherits the memory-intensive characteristic of poly-  
251 nomial computations. This highlights the reality that  
252 greater computational power is not the primary impetus  
253 for accelerating FHE schemes. Instead, the predominant  
254 constraint that curtails the performance is the efficiency  
255 of data movement. To further quantify the impact of  
256 bandwidth on performance and to gain insight into the  
257 memory/compute characteristics of FHE operations, we  
258 build a roofline model [27] using the NVIDIA GeForce  
259 RTX 3090 as an illustrative instance to show the primary  
260 bottleneck in FHE operations. We adopt FHE encryption  
261 parameters as in [13] (e.g.,  $N = 2^{16}$ ,  $L = 45$ ) for the  
262 illustrative example. As shown in Fig. 2, the horizontal  
263 roof represents the peak computation capacity (i.e., 35.6  
264 TFLOPS) while the diagonal roof expresses the memory  
265 bandwidth (i.e., 936 GB/s) of the target hardware plat-  
266 form. The x-axis represents AI, which is a ratio of  
267 computation to memory access during the execution  
268 progress of CKKS operations. Overlaying the CKKS  
269 APIs on this model reveals a significant insight: all



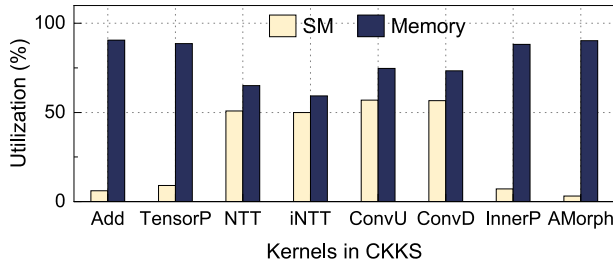


Fig. 3. SM and memory utilization during different CKKS kernels runtime.

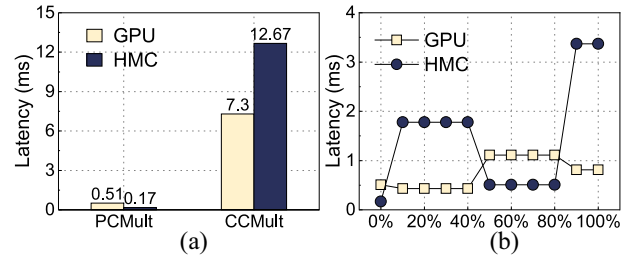


Fig. 4. Different affinity tendencies of different CKKS operations, refer to Section V for detailed settings. (a) CKKS operation. (b) Breakdown analysis.

270 CKKS APIs exhibit significantly low AI. In particular,  
 271 these operations are situated below the diagonal roof  
 272 and are far from fully exploiting the peak computational  
 273 power of the hardware platform. This finding suggests  
 274 that the efficiency of data movement is constraining the  
 275 performance of these APIs, and simply deploying them  
 276 to more powerful hardware devices may result in little  
 277 boost in the overall performance.

278 2) *Observation 2—Different CKKS Basic Kernels Show*  
 279 *Varying Memory Bandwidth and Computational*  
 280 *Requirements:* Our evaluation results also reveal that  
 281 although all CKKS APIs exhibit memory-intensive  
 282 characteristics, the underlying basic kernels within  
 283 CKKS APIs prioritize computing power and bandwidth  
 284 requirements differently. We ran all CKKS kernels on  
 285 an NVIDIA Geforce 3090 platform and accessed their  
 286 resource utilization using the Nsight Compute [28]  
 287 tool. The SM and Memory utilization results shown  
 288 in Fig. 3 demonstrate while nearly all kernels exhibit  
 289 memory-intensive behavior, certain kernels (e.g., *NTT*,  
 290 *Conv*) also impose substantial demands on computation  
 291 resources. This suggests that there might be kernels more  
 292 suited for execution on the GPU engines. To further  
 293 investigate the aforementioned issues, we choose two  
 294 of the most representative CKKS APIs (i.e., *PCMult*: a)  
 295 plaintext-ciphertext multiplication, comprising *TensorP*  
 296 and b) *CCMult*: Ciphertext-ciphertext multiplication,  
 297 comprising a series of CKKS basic kernels) and  
 298 run them on GPU and HMC processing engines  
 299 (for details, see Section V), respectively. Fig. 4(a)  
 300 illustrates the execution latency of these two APIs on  
 301 different platforms, whereas Fig. 4(b) provides the time  
 302 consumption of different stages during the execution of  
 303 *CCMult* on GPU and HMC sides. The results shown in  
 304 Fig. 4(a) and (b) reveal two critical phenomena:

- 305 a) Different CKKS APIs exhibit distinct preferences  
 306 for HMC sides. Specifically, *PCMult* tends HMC  
 307 side, while *CCMult* demonstrates a preference for  
 308 the GPU side.  
 309 b) Even within a single CKKS API, the various stages  
 310 of its internal execution process display differing  
 311 affinities for HMC side. This phenomenon primarily  
 312 stems from the fact that *CCMult* encompasses  
 313 all CKKS kernels outlined in Section II-A except  
 314 for *AMorph*, and these kernels possess varied  
 315 affinities with the HMC side.

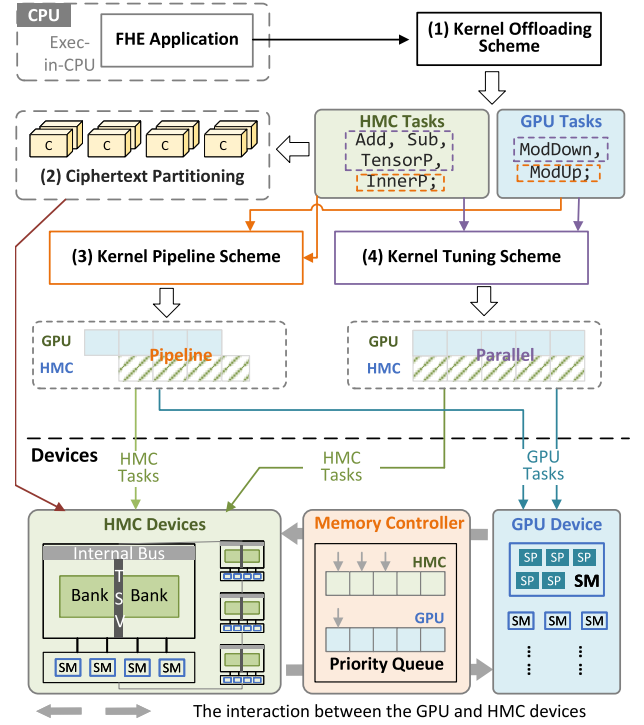


Fig. 5. Overview of HMC-FHE.

316 These two observations reveal the importance of bandwidth  
 317 in the efficiency of CKKS operations. Additionally, they high-  
 318 light another crucial point: blindly offloading FHE operations  
 319 to HMC processing engines is not always advantageous. This  
 320 is because certain FHE kernels still require a delicate balance  
 321 between computational and bandwidth requirements. In this  
 322 study, we aim to develop a kernel-level scheduling framework  
 323 leveraging the GPU-HMC heterogeneous architecture to effi-  
 324 ciently accelerate FHE operations.

### III. HMC-FHE FRAMEWORK

325 Fig. 5 provides an overview of the proposed HMC-FHE  
 326 design. A typical GPU-HMC architecture is adopted as the  
 327 accelerator hardware [17], [18]. This architecture integrates  
 328 three distinct types of devices, each with specialized capabilities:  
 329 1) GPU, serving as the primary processing engine, focuses  
 330 on computation-intensive tasks; 2) HMC, serving as a sec-  
 331 ondary processing engine and the global memory for the GPU,  
 332 focuses on memory-intensive tasks; and 3) CPU, serving as an  
 333 auxiliary processing device, receives applications and work-  
 334 loads and plays a pivotal role in overseeing task scheduling  
 335

and resource management between GPU and HMC engines. To eliminate redundant data transfers between GPU and HMC devices, the HMC devices are directly connected to the GPU side via their memory links on a silicon interposer, replacing the traditional GPU global memory interface. Additionally, the interaction and data communication among these three devices is facilitated through the packet-based protocol (as mentioned in Section II-B), which is managed by the HMC devices, ensuring reliable and efficient data transfer within the system.

To accelerate an FHE application, we propose the HMC-FHE framework to automatically map and schedule the FHE operations on the GPU-HMC architecture, comprising four crucial components: 1) a fine-grained kernel offloading scheme; 2) a ciphertext partitioning scheme; 3) a kernel execution pipeline scheme; and 4) a kernel tuning scheme. These components collaborate to address three challenges as mentioned in Section I.

The fine-grained kernel offloading scheme is utilized to distribute various CKKS basic kernels to GPU or HMC engines based on their resource affinity. To enhance the parallelism of CKKS basic kernels assigned to HMC engines, we propose a ciphertext partitioning scheme aimed at achieving data parallelism across multiple HMC devices, which is achieved by mapping each ciphertext to multiple HMC devices. The kernel execution pipeline scheme leverages the loose data dependencies within the data flow graph (DFG) of those compound CKKS operations and allows their composed kernels to be executed between GPU and HMC engines in a pipeline fashion. Finally, to fully harness the capabilities of both the GPU and HMC engines, a kernel tuning module is employed to dynamically reallocate kernels between the GPU and HMC engines based on runtime status. This ensures the maximization of system parallelism.

#### IV. HMC-FHE COMPONENT DESIGN

In this section, we will provide the design details of the proposed four key functional components of HMC-FHE.

##### A. FHE Kernel Offloading Model

Based on the analysis of FHE operations as shown in Section II-C, we propose a fine-grained FHE kernel offloading model to offload these underlying CKKS basic kernels smartly between GPU and HMC processing engines. The fundamental principle lies in the AI value of a kernel, which notably impacts its ideal processing location: kernels with lower AI tend to favor the HMC processing engine, as they require less computational power. Conversely, kernels with higher AI, signaling a greater demand for computational resources, are more effectively to be offloaded to the GPU. This concept is further illustrated in Fig. 6, showcasing a roofline model example utilizing an NVIDIA GPU with multiple HMC cubes (for details, see Section V-A). Fig. 6 show that kernels with AI values below the Critical AI (CAI, defined as the ratio of the HMC peak computation capacity to GPU bandwidth) can effectively utilize the computing resources available at the HMC processing engine. However, when the AI of a kernel

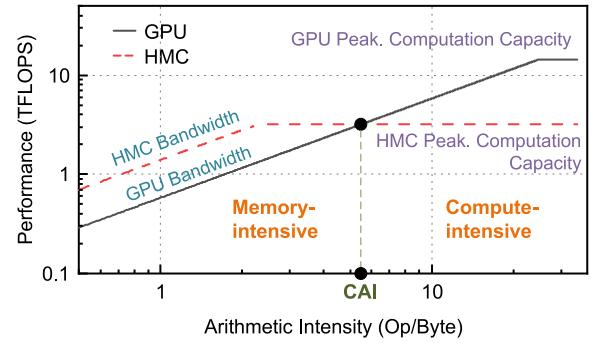


Fig. 6. Roofline model based on the GPU-HMC architecture adopted in this article.

surpasses the CAI, executing it on the GPU processing engine becomes more beneficial.

However, obtaining the AI of various CKKS kernels with different parameters in advance presents a challenge. To address this, we formulate a linear programming model for these CKKS kernels to quantify their computation and memory access ratio. This linear programming model comprises three key metrics, namely, global memory AI (MAI, the value is denoted as  $m$ ), shared memory AI (SMAI, the value is denoted as  $s$ ), and the ratio of shared memory instruction (SMI), count to the total instruction (TI) count (SMI: TI ratio, the value is denoted as  $r$ ). MAI describes the computation-to-data access ratio when data resides in memory from an instructional perspective, while SMAI represents the same ratio when data resides in the cache. Both metrics function akin to AI, with lower values indicating reduced computational demands. Additionally, SMI:TI delineates the ratio of instructions accessing cached data to TIs, serving as a weighting mechanism to gauge the interplay between MAI, SMAI, and AI. A higher value suggests that SMAI exerts a stronger influence on AI. Note that all metrics are static and can be readily obtainable through source code parsing. Consequently, the host can seamlessly execute the kernel offloading process in real-time with negligible overhead

$$\text{Affinity} = \begin{cases} \text{HMC}, & \text{if } x_{1-r}x_m + c_1x_r x_s < \text{CAI} \\ \text{GPU}, & \text{otherwise.} \end{cases} \quad (1)$$

Hence, the kernel offloading model is as shown in (1), where  $c_1$  is a constant, representing the linear factors of the mapping of SMAI to CAI.  $x_m$ ,  $x_s$ ,  $x_r$  represent the MAI, SMAI, SMI:TI Ratio of kernel  $x$ , respectively. Accordingly, if the output of the model is less than CAI, it indicates that the current kernel is memory-intensive, implying a higher priority for its bandwidth requirement. Consequently, the kernel will be offloaded to the HMC side. Otherwise, the kernel will be offloaded to the GPU side. Table II illustrates the outcomes of various metrics and the affinity of different CKKS kernels within the context described in Section II-A.

##### B. Bandwidth-Aware Ciphertext Partitioning

To leverage the parallelism offered by the multiple HMCs, those basic kernels offloaded to the HMC can be processed in two ways: one approach is assigning a series of basic kernels to different HMC engines by a proper schedule mechanism. Each

TABLE II

METRICS AND THE AFFINITY RESULTS UNDER THE SETTINGS MENTIONED IN SECTION V-A. LEGEND: I) MAI AND SMAI (MEMORY TO COMPUTE RATIO =  $0 \leq L < 1 \leq M < 5 \leq H$ ) AND II) RATIO (SHARED INS. TO INS. =  $0 \leq L < 0.1 \leq M < 0.5 \leq H$ )

Kernel	MAI	SMAI	SMI:TI Ratio	Affinity
Add	L	L	L	HMC
Sub	L	L	L	HMC
TensorP	L	L	L	HMC
iNTT	H	M	M	GPU
NTT	H	M	M	GPU
ConvU	H	M	M	GPU
ConvD	H	M	M	GPU
InnerP	L	L	L	HMC
AMorph	L	L	L	HMC

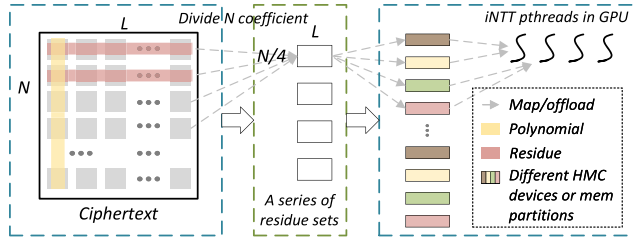


Fig. 7. Overview of Ciphertext partitioning scheme and assume that the number of HMC devices is 4. The ciphertext is composed of  $L$  polynomials and each polynomial contains  $N$  coefficients. The coefficients from the same position of the  $L$  polynomials are denoted as residue.

HMC engine manages a complete ciphertext independently, treating each HMC as an autonomous entity. Alternatively, one can accelerate a single basic kernel using multiple HMCs, with each HMC handling a portion of a single ciphertext. In this setup, all HMC engines function collectively as a unified entity. The former is heavily dependent on the DFG of applications. If there are data dependencies between basic kernels, some HMC processing engines may remain idle, leading to inefficient resource utilization. The latter allows multiple HMC engines to collaboratively execute a CKKS basic kernel, effectively guaranteeing resource utilization regardless of the DFG of applications. In this article, our focus lies on this latter method.

As mentioned in Section II-B, it was noted that the transfer latency between HMC devices exceeds three times that of local access within an HMC. Hence, the ciphertext partitioning scheme must be carefully designed to ensure its performance gain. We introduce a bandwidth-aware ciphertext partitioning (BaCP) scheme, which fulfills the following three criteria by maintaining the coherence of parallel and data access patterns within the CKKS basic kernel: 1) ensuring the parallel processing capabilities of multiple HMC engines; 2) attaining load balance across the multiple HMC engines; and 3) minimizing data transfer between HMC devices.

As shown on the left side of Fig. 7, the original dense ciphertext is organized as an  $N * L$  matrix. This structure can be understood as comprising either  $L$  polynomials (each of degree  $N$ ) or  $N$  residues (each of degree  $L$ ). In the BaCP scheme, the ciphertext is first decomposed into four parts, with each part containing  $(N/4)$  residues (assuming there are four HMC devices). Subsequently, each part is evenly distributed

across different HMC devices, as illustrated on the right side of Fig. 7. For instance, considering HMC one, it manages residues indexed  $1_{st}, 4_{th}, 8_{th}, \dots, (N-3)_{th}$ .

This partitioning scheme brings two advantages.

- 1) *Parallelism and Load Balancing*: The CKKS basic kernel, initially processed by a single HMC engine, is transformed into four parallel subkernels of equal scale and function. These four subkernels can then be concurrently processed by four HMCs, facilitating parallel execution.
- 2) *Decoupling of Ciphertext Structure and Access Patterns*: By decoupling the inherent structure of  $L$  polynomials, our approach enables the alignment of data access patterns and parallel patterns within CKKS basic kernels, fostering collaboration between them. Consider the scenario where the HMC side processes the *InnerP* kernel, which involves data movement across  $L$  polynomials and conducts multiply-accumulation (MAC) operations. The BaCP scheme redefines it as four smaller subkernels, each handling inputs of size  $(N/4) * L$ . Each subkernel then executes  $(N/4)$  independent MAC operations aimed at  $(N/4) * L$  residues. Within this setup, every MAC operation utilizes elements from the same residue, ensuring that all data access occurs locally within the vaults of a single HMC device.

The BaCP scheme aligns well with CKKS kernels offloaded to the HMC engines, except for *AMorph*. For kernels, such as *Add*, *Sub*, and *TensorP*, which involve element-wise operations like multiplication or addition, the BaCP scheme avoids conflicts between parallel modes and data access patterns. However, *AMorph*, due to its data access pattern being the opposite of *InnerP*, necessitates some level of data transfer between HMCs under the BaCP scheme. Nonetheless, considering the infrequent calls of *AMorph* and the relatively minimal data transfer compared to *InnerP*, the overall efficiency of the BaCP scheme is minimally impacted.

### C. Kernel Execution Pipeline

While the offloading mechanism and ciphertext partitioning scheme improve the execution efficiency of CKKS kernels, the inherent data dependencies among these kernels necessitate their execution on the GPU and HMC engines in a serial manner. This introduces idle periods that reduce system parallelism and throughput. Taking the *CCMult* API as an example, which comprises four sequential CKKS operations (*TensorP*  $\rightarrow$  *ModUp*  $\rightarrow$  *InnerP*  $\rightarrow$  *ModDown*), the kernel offloading model assigns *TensorP* and *InnerP* to the HMC engines (as indicated in Table II). Consequently, the processing engines within both the GPU and HMC engines will alternate between idle states, as depicted in Fig. 8(a). To tackle this issue, our objective is to identify loose data dependencies from interkernel within the DFG and adjust the execution sequence of certain kernels accordingly. This allows for the asynchronous execution and pipeline processing of different CKKS kernels across the GPU and HMC sides, thereby reducing the idle period of the system.

Our main focus lies on CKKS compound operations, particularly on compound kernels (e.g., *Bootstrapping*) or CKKS



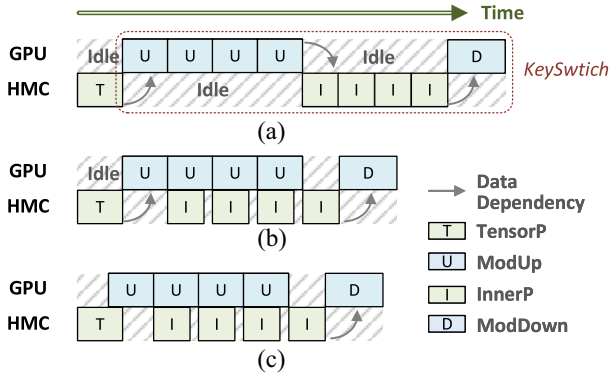


Fig. 8. Pipelined execution inside GPU-HMC architecture. *ModUP* (*ModDown*) consists of *NTT*, *iNTT*, and *ConvU* (*ConvD*). The parameter  $dnum = 4$ . (a) Execution flow of a *CCMult*. (b) Loose data dependency 1. (c) Loose data dependency 2.

518 APIs (e.g., *CCMult*, *Rotate*) that incorporate the *KeySwitch*  
 519 kernel. These compound operations, including *KeySwitch*, are  
 520 primarily focused on for the following reasons.

- 521 1) They are the most time-consuming operations. A sig-  
 522 nificant proportion of time in applications is consumed  
 523 by a large number of *KeySwitch* calls, constituting over  
 524 90% of the time cost of applications [4], [29].
- 525 2) According to the kernel offloading model, these com-  
 526 pound operations uniquely exhibit the characteristic  
 527 of alternating between the GPU and HMC sides. In  
 528 contrast, basic kernels like *Add* and *Sub* do not support  
 529 pipelining across both sides. After a comprehensive  
 530 analysis of the DFG for the mentioned operations, we  
 531 identify two primary types of loose data dependencies.  
 532 To illustrate these differences, we utilize the *CCMult*  
 533 API as a case study:

534 a) *Loose Data Dependency 1*: For the data depen-  
 535 dency within *KeySwitch* (*ModUp*  $\rightarrow$  *InnerP*  $\rightarrow$   
 536 *ModDown*), we employ the General *KeySwitch*  
 537 scheme [30] to divide the *ModUp*  $\rightarrow$  *InnerP*  
 538 sequence into  $dnum$  smaller stages. At this point,  
 539 the output of each *ModUp* serves as the input for  
 540 the corresponding *InnerP*, and the  $dnum$  *ModUp*  
 541  $\rightarrow$  *InnerP* stages operate in complete parallelism.  
 542 Consequently, we can utilize a pipelining approach  
 543 to alternately execute *ModUp* and *InnerP* on the  
 544 GPU and HMC engines, thereby implicitly conceal-  
 545 ing resource idleness. Subsequently, the data  
 546 flow of *CCMult* transitions from the configuration  
 547 depicted in Fig. 8(a) to that illustrated in Fig. 8(b).

548 b) *Loose Data Dependency 2*: Another instance of  
 549 loose data dependence occurs at the intersection  
 550 of the *TensorP* kernel and the *KeySwitch* kernel.  
 551 Specifically, *TensorP* computes four element-wise  
 552 multiplications and generates three outputs, one  
 553 of which serves as the input to *KeySwitch*.  
 554 This implies that we can rearrange the execu-  
 555 tion sequence between *TensorP* and *KeySwitch*,  
 556 allowing *TensorP* to first compute the *KeySwitch*  
 557 input. Subsequently, the data dependency between  
 558 *TensorP* and *KeySwitch* is eliminated, enabling

559 their parallel execution at runtime. After restructur-  
 560 ing the data flow, the data flow of *CCMult*  
 561 transitions from the configuration depicted in  
 562 Fig. 8(b) to that illustrated in Fig. 8(c), further  
 563 enhancing system efficiency.

564 Note that the aforementioned loose data dependencies are  
 565 general and can be extended to other operations that include  
 566 *KeySwitch*. Therefore, the two pipeline schemes described also  
 567 apply to *Rotate* and *Bootstrapping*.

#### D. Kernel Tuning

568 The kernel execution pipeline scheme primarily guarantees  
 569 parallel processing capability in the system for specific kernels.  
 570 However, in other scenarios, such as when the GPU engine  
 571 executes *ModDown* or the HMC engine handles basic kernels  
 572 like *Add* or *Sub*, the processing engines on the opposite side  
 573 remain idle.

574 To address this issue, we propose a kernel tuning scheme  
 575 aimed at optimizing resource utilization throughout the entire  
 576 runtime. This scheme entails extracting subkernels from the  
 577 active GPU or HMC engines and reallocating them to idle  
 578 HMC or GPU engines. This ensures a more balanced and  
 579 efficient utilization of resources across the architecture. It  
 580 is important to note that this approach deviates from the  
 581 kernel offloading model, but the enhanced utilization of system  
 582 resources yields performance gains.

583 To ensure the efficiency of kernel execution, our kernel  
 584 tuning scheme adheres to two design principles: 1) ensuring  
 585 load balancing between the GPU and HMC engines and  
 586 2) minimizing the data transfer overhead between multiple  
 587 HMC devices, particularly when extracting subkernels from  
 588 the GPU engine and assigning them to the HMC engines. Next,  
 589 we provide a detailed explanation of the kernel tuning scheme  
 590 from both the GPU and HMC perspectives.

591 *Kernels in HMC*. For basic kernels offloaded to the HMC  
 592 engines (e.g., *Add*, *Sub*), since they lack specific layout con-  
 593 straints for ciphertext elements (as discussed in Section IV-B),  
 594 it is feasible to extract subkernels and reoffload them to  
 595 the GPU engine. Taking the tuning process of the *Add*  
 596 kernel illustrated in Fig. 9(a) as an example, its ciphertext is  
 597 partitioned into two parts, with the GPU and HMC engines  
 598 handling one portion each, respectively. Moreover, to ensure  
 599 load balance, the ratio of data scale processed by the GPU and  
 600 HMC engines adheres to

$$\frac{\text{GPU}_{\text{workload}}}{\text{HMC}_{\text{workload}}} = \frac{\text{Roofline}_G(\theta(x))}{\text{Roofline}_H(\theta(x))} \quad (2)$$

603 where  $\theta(x) = x_{1-r}x_m + c_1x_r x_s$  (as discussed in Section IV-A),  
 604 represents the AI of kernel  $x$ .  $\text{Roofline}_G(\theta(x))$  and  
 605  $\text{Roofline}_H(\theta(x))$  denote the available peak performance of  
 606 kernel  $x$  on the roofline models of GPU and HMC,  
 607 respectively. Therefore, the ratio of data scale (i.e.,  
 608  $\text{GPU}_{\text{workload}}/\text{HMC}_{\text{workload}}$ ) processed by both sides aligns  
 609 with the ratio of their performance for associated kernels.

610 *Kernels in GPU*: For *ModDown*, the situation is consider-  
 611 ably more complex. This complexity arises primarily because  
 612 the *NTT* and *iNTT* kernels within *ModDown* necessitate data  
 613 permutation across coefficients of each polynomial, whereas

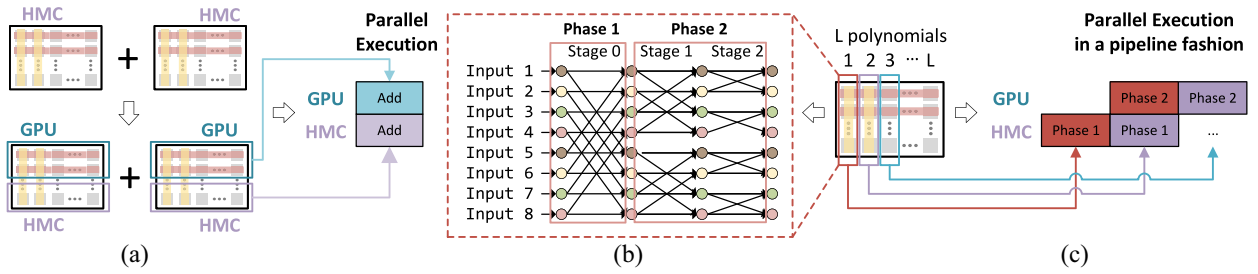


Fig. 9. Kernel tuning scheme. (a) and (c) describes the kernel tuning details between GPU and HMC engines; (b) illustrates the execution progress of an 8-input *NTT*, where different colors represent elements located in different HMC devices.

614 the ciphertext partitioning scheme does not allocate a complete  
 615 polynomial to an HMC device. Consequently, when these  
 616 kernels are offloaded to the HMC engines, it initiates extensive  
 617 data exchanges between HMC devices, thereby diminishing  
 618 the data movement efficiency of the system. To address this  
 619 issue, the kernel tuning scheme adopts a 2-phase *NTT* method  
 620 based on the relationship between the data access pattern  
 621 within the *NTT* and the BaCP scheme.

622 Fig. 9(b) outlines the execution details of an 8-input *NTT*,  
 623 wherein the computation process is segmented into three  
 624 stages, each comprising a series of butterfly operations [31].  
 625 Suppose this kernel is reallocated to HMC engines, with the  
 626 polynomial interleaved across the four HMC devices following  
 627 the BaCP approach. In Stage 0, input  $x$  undergoes butterfly  
 628 operations with input  $x + 4$  ( $1 \leq x \leq 4$ , stride = 4), where  
 629 the inputs for each butterfly operation originate from the same  
 630 HMC device. However, in Stage 1 or 2, the stride for the  
 631 butterfly operations is 2 or 1, respectively, necessitating data  
 632 exchanges between HMC devices.

633 We can extend this scenario to larger *NTT* inputs. For any  
 634 given  $N$ , any HMC device can execute the  $0 \sim \log_2(N/4) - 1$   
 635 stages of *NTT* without data transfer between HMC devices.  
 636 This observation motivates us to split the *NTT* kernel into  
 637 two phases and allocate phase 1 to the HMC engines.  
 638 Consequently, the HMC engines can perform a portion of the  
 639 *NTT* calculation under the BaCP scheme.

640 As depicted in Fig. 9(c), the GPU and HMC engines  
 641 sequentially execute a 2-phase *NTT* for a single poly-  
 642 nomial. For  $L$  polynomials within a ciphertext, which can be  
 643 processed in parallel, the GPU and HMC effectively utilize  
 644 system resources without idleness through pipeline execution.  
 645 Additionally, to sustain load balance, the number of stages in  
 646 phase 1 and phase 2 also conforms to (2).

647 Note that this approach can be readily applied to the *INTT*  
 648 stage with minimal modifications. Since the *INTT* kernel is  
 649 essentially the inverse of the *NTT*, the processing sequence  
 650 between the GPU and HMC is simply inverted. In other words,  
 651 the HMC handles phase 2 of *INTT*, while the GPU undertakes  
 652 phase 1 of *INTT*.

## 553 V. EVALUATION

654 In this section, we will provide a comprehensive overview  
 655 of the experimental setup, workloads employed, evaluation  
 656 results, and comparison with other implementations.

TABLE III  
 PLATFORM SPECIFICATIONS

Baseline Configuration		
GPU platform	Total 56 SMs, 32 SIMT, GTO warp schedule, 1290MHz, L1 128 KB (per SM), Main Memory DDR6 (480 GB/s)	
GPU-HMC Configuration		
GPU side	Total 40 SMs, 32 SIMT, GTO warp schedule, 1290MHz, L1 128 KB (per SM)	
HMC side	Device	4 HMC cubes, total 32 GB capacity, 312.5MHz, 320 GB/s Internal bandwidth
	Processing logic	Total 16 SMs (4 SMs per HMC cube), 32 SIMT, 1290MHz, L1 16 KB (per SM)
	Crossbar	6-cycle all-to-all crossbar
	Serials links	10-cycle latency, including 3.2 ns for SerDes, 2 pJ/bit (1.5 pJ/bit/cycle when it is idle)
	3D DRAM timings	$t_{CK} = 1.6$ , $t_{CAS} = 11.2$ , $t_{RD} = 11.2$ , $t_{RAS} = 22.4$ , $t_{RP} = 11.2$ , $t_{WR} = 14.4$

### 557 A. Environment Setup

558 *Experimental Setup:* We implement the prototype of HMC-  
 559 FHE by integrating the latest GPGPU-Sim v4.0 [32] with the  
 560 HMC simulator CasHMC [33]. Table III provides the details  
 561 of the platform specifications. We configure the GPU side  
 562 with 40 SM cores and each HMC processing engine with  
 563 four SM cores (< 30 W, completely satisfying the thermal  
 564 feasibility [34]) underneath a 3-D memory stack. We configure  
 565 the platform with four HMC cubes and each cube of the  
 566 HMC is embedded in a memory partition, which is directly  
 567 connected to the weakened SMs located in the logic layer  
 568 through an interconnection network. The external transmission  
 569 latency of the memory partition is configured using the link  
 570 model in BookSim [35]. In addition, we use the AccelWatch  
 571 model [36] to evaluate EE under different schemes, where the  
 572 energy parameters follow [37]. The diagram of the integrated  
 573 architecture of GPU and HMC is as in Fig. 1.

574 *Simulator Validation:* The GPU-HMC architecture, essential  
 575 to our simulations, has been extensively used in prior research,  
 576 setting a reliable precedent for its use in academic studies.  
 577 This guarantees that our architectural assumptions are based  
 578 on tested and proven models. Our simulations are carried  
 579 out using GPGPU-Sim and CasHMC, both of which are  
 580 tools well-recognized within the research community for their  
 581 accuracy and reliability. Moreover, previous research [32] has  
 582 shown that GPGPU-Sim offers over 85% accuracy in its  
 583 simulation results. This high level of precision significantly  
 584 boosts the credibility of our simulation outcomes, indicating  
 585 that they are a credible reflection of real-world performance.



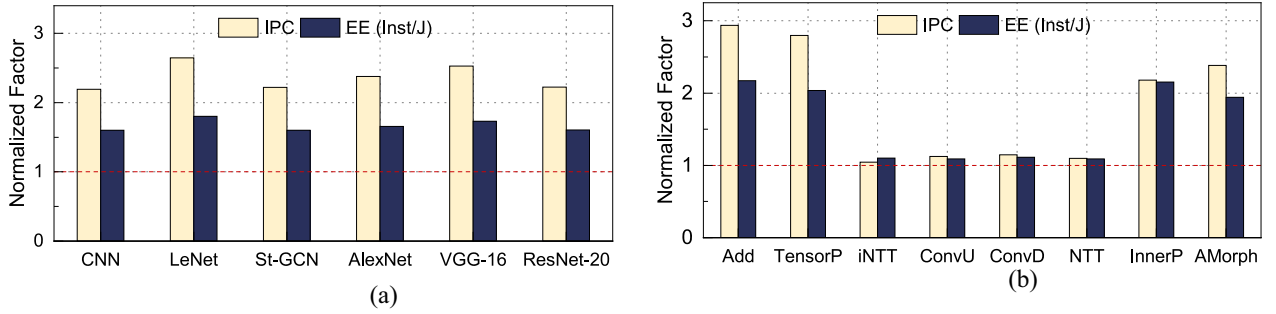


Fig. 10. (a) Throughput improvement and energy saving. (b) Performance improvement for different CKKS basic kernels.

**Baseline:** We compare this design with a conventional GPU architecture that has 56 SM cores as the baseline. Note that our comparison is based on a fairground of similar computational power and bandwidth resources (i.e., both of them have the same peak computing power and external interconnection bandwidth) We also evaluate this framework with more SM cores, see the details in Section V-D.

**Benchmark:** Six neural network models, CNN [4], LeNet [38] with dataset MNIST, St-GCN [2] with dataset NTU-RGB+D, AlexNet [38], VGG-16 [38], and ResNet-20 [16] with dataset CIFAR-10, that are typically adopted in privacy protection machine learning are used for performance evaluation. Among them, the CNN, LeNet, and AlexNet models are used for evaluating the cases without *bootstrapping*, and the other models are used for evaluating the cases with *bootstrapping*. As a result of the distinct multiplicative depth requirements (i.e., without/with *bootstrapping*) for these two types of models, we employ two separate sets of CKKS parameters denoted as  $(N = 2^{16}, L = 45)$  and  $(N = 2^{14}, L = 20)$ , respectively. The security level  $\lambda > 98$ .

**Implementation:** For the implementation of CKKS, we adopt `ckks-gpu-core` [13] (the state-of-the-art GPU implementation of the CKKS scheme). We adapt parts of the implementation to enable its execution on HMC-FHE. During the runtime phase, ciphertexts along with associated constant parameters (e.g., twiddle factors in  $(I)NTT$ , RNS modulus generated by the CRT) involved in CKKS operations are generated in host memory and then efficiently distributed to multiple HMC devices using BaCP methodology. To effectively manage the substantial number of evaluation keys (*evks*) during model execution tasks, particularly when invoking *Bootstrap*, we devise a memory pool system. This system oversees memory allocations tailored to HMC, mitigating the risk of memory overflow. Under this framework, all *evks* are initially stored on the host side for ease of management and access. When these keys are needed for computation, the system copies an instance of *Evk* from the host side, formats it into a format BaCP suitable for parallel data processing, and sends it to multiple HMC devices.

## B. Overall Performance

To assess the efficacy of the acceleration framework, we perform a comprehensive evaluation of end-to-end inference performance using the benchmarks mentioned above. The obtained results are subsequently compared with the baseline

(i.e., GPU platform), with instructions per cycle (IPC) and EE being the primary metrics of comparison.

Fig. 10(a) provides the performance speedup and energy savings achieved across all workloads. Regarding inference latency, the HMC processing engines effectively alleviate the “memory wall” bottleneck within the CKKS kernels. Consequently, CNN, LeNet, St-GCN, AlexNet, VGG-16, and ResNet-20 exhibit performance improvements of  $2.19\times$ ,  $2.65\times$ ,  $2.22\times$ ,  $2.38\times$ ,  $2.53\times$ ,  $2.22\times$  compared to the baseline, respectively. Moreover, due to the reduced power overhead of data movement, our acceleration framework decreases energy consumption across different benchmarks by 59.95%, 80.34%, 60.22%, 65.22%, 73.32%, and 60.34%, respectively. Moreover, since all workloads show similar performance gains, HMC-FHE can be well-compatible with different CKKS parameters.

To provide additional evidence of the acceleration benefits conferred by HMC-FHE across various CKKS kernels, we monitored the execution progress of the ResNet-20 models. We recorded the average execution time of different kernels and evaluated their corresponding EE from a fundamental kernel perspective. As shown in Fig. 10(b), as expected, the kernels (*Add*, *TensorP*, *InnerP*, *AMorph*) exhibit substantial speedup, with improvements of  $2.93\times$ ,  $2.79\times$ ,  $2.18\times$ , and  $2.38\times$ , respectively. Conversely, the kernels (*iNTT*, *NTT*, *ConvU*, *ConvD*) demonstrate comparatively modest speed increases of 4.53%, 12.39%, 14.46%, and 9.68%, respectively. This is as expected. As shown in Table II, the kernels (*Add*, *TensorP*, *InnerP*, *AMorph*) are offloaded to HMC processing engines. Hence, they can benefit from the larger internal bandwidth, and lower latency provided by the HMC devices, resulting in significant performance improvement. On the contrary, the kernels (*iNTT*, *NTT*, *ConvU*, *ConvD*) are performed in the GPU engine, the incompatible memory resource usage makes the acceleration effect of these kernels less obvious. Note that although various basic kernels experience different levels of speedup, their individual performance does not directly influence the overall application speedup at a global level. With the pipelined design and kernel tuning scheme, the two types of kernels running on GPU and HMC sides are executed in a parallel fashion, which to some extent conceals a significant portion of time overhead.

## C. Breakdown Analysis

To evaluate the benefit of each technique toward overall performance improvement, we form combinations of different

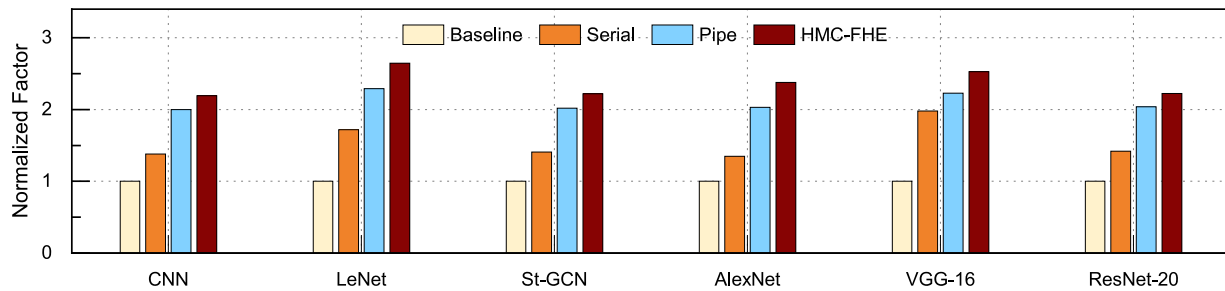


Fig. 11. Breakdown analysis of different schemes across all workloads.

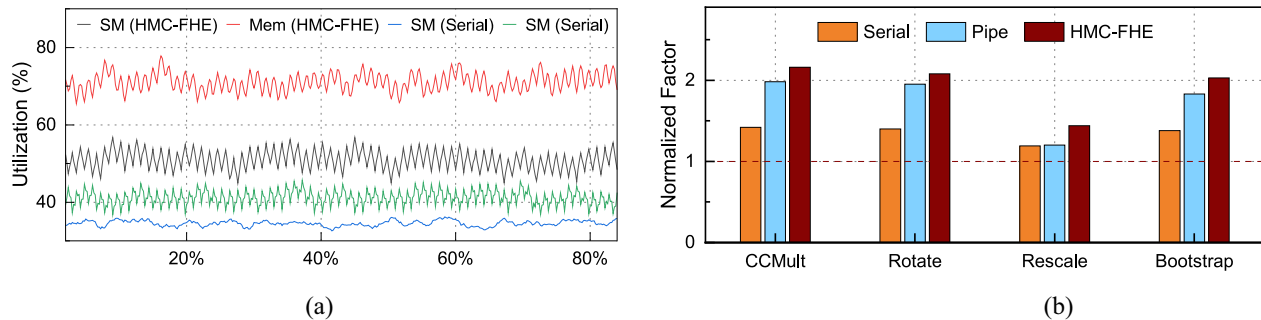


Fig. 12. (a) Resource utilization under ResNet-20. (b) Breakdown analysis of CKKS compound operations.

774 schemes to elucidate the advantages of these combined tech- 808  
 775 niques both at the application and primitive levels. The specific 809  
 776 combinations of schemes are outlined below. 810

- 777 1) *Baseline*: GPU only implementation. 811
- 778 2) *Serial*: Only implement the kernel offloading 812  
 779 (Section IV-A) and the Ciphertext partition scheme 813  
 780 (Section IV-A), with all CKKS primitives serially 814  
 781 executed on HMC-FHE. 815
- 782 3) *Pipe*: Supplement the FHE operation scheduler 816  
 783 (Section IV-C) to support pipelined execution based 817  
 784 on Serial, but without implementing the kernel tuning 818  
 785 scheme (Section IV-D). 819
- 786 4) *HMC-FHE*: Integrating all the technical schemes. 820

787 1) *Application-Level*: Fig. 11 presents the end-to-end 821  
 788 inference performance of all six models with different 822  
 789 technique combinations. By solely employing the *Serial* 823  
 790 scheme, we observe an enhancement in the inference 824  
 791 performance of these NN models compared to the *Baseline* by 825  
 792  $1.38\times$ ,  $1.72\times$ ,  $1.41\times$ ,  $1.35\times$ ,  $1.98\times$ , and  $1.42\times$ , respectively. 826  
 793 This improvement is primarily credited to the advantages 827  
 794 facilitated by HMC devices. Furthermore, the varied gains 828  
 795 observed across models stem from the differing proportions 829  
 796 of FHE operations conducted during their execution. For 830  
 797 instance, VGG-16 is predominantly characterized by *Add* and 831  
 798 *TensorP* kernels, with relatively fewer calls to *KeySwitch*. 832  
 799 As a result, the *Serial* solution can fully exploit the 833  
 800 benefits of HMC in this scenario. In comparison to the 834  
 801 *Serial* solution, *Pipe* yielded an enhancement in inference 835  
 802 performance by  $1.44\times$ ,  $1.33\times$ ,  $1.43\times$ ,  $1.51\times$ ,  $1.13\times$ , and 836  
 803  $1.44\times$ , respectively. This improvement primarily arises from 837  
 804 the optimization of *KeySwitch* operations, which augments 838  
 805 system parallelism while fully utilizing the architecture's 839  
 806 overall resources. Lastly, HMC-FHE incorporates the Kernel 840  
 807 tuning scheme atop *Pipe*, further refining the overall resource

utilization of the architecture and extending parallel processing 808  
 between the GPU and HMC sides to all phases of model 809  
 inference. Consequently, HMC-FHE achieved a performance 810  
 improvement of  $1.10\times$ ,  $1.16\times$ ,  $1.10\times$ ,  $1.17\times$ ,  $1.13\times$ , and 811  
 $1.09\times$  compared to *Pipe*. 812

To illustrate the variations in resource utilization across 813  
 different technology combinations, using ResNet-20 as an 814  
 example, Fig. 12(a) displays the average resource utilization 815  
 (i.e., SM, Memory) during the execution process under both 816  
 the *Serial* and *HMC-FHE* schemes. By employing kernel 817  
 pipeline execution and kernel tuning, parallelism is optimized 818  
 on GPU and HMC, reducing idle stalls. Consequently, the 819  
 SM (weighted average of SM utilization on both GPU and 820  
 HMC sides) and memory utilization in the *HMC-FHE* scheme 821  
 exhibit a significant increase compared to the *Serial* scheme. 822

2) *Compound Operations-Level*: Fig. 12(b) depicts the 823  
 performance improvement of all compound CKKS APIs 824  
 in the aforementioned models under different technology 825  
 combinations. Similar to the *application-level* scenario, the 826  
*Serial* scheme exhibits varying acceleration effects for differ- 827  
 ent CKKS compound operations, yielding improvements of 828  
 $1.42\times$ ,  $1.40\times$ ,  $1.19\times$ , and  $1.38\times$ , respectively. Furthermore, 829  
 with the additional deployment of kernel pipeline execution, 830  
 the performance of these compound operations is further 831  
 enhanced by  $1.39\times$ ,  $1.39\times$ ,  $1.00\times$ , and  $1.33\times$ , respectively. 832  
 In this context, since the *Rescale* does not integrate the 833  
*KeySwitch* kernel internally, the kernel pipeline scheme is 834  
 ineffective for it. Finally, with the Kernel tuning scheme, the 835  
 performance sees  $1.10\times$ ,  $1.07\times$ ,  $1.21\times$ , and  $1.11\times$  speedup. 836

#### D. Scalability 837

One of the advantages of HMC-FHE is its exceptional scala- 838  
 bility. As mentioned above, the ciphertext partitioning scheme 839

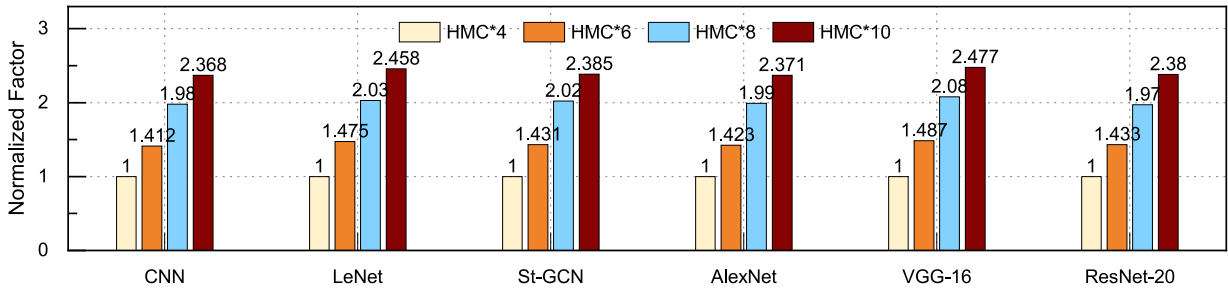


Fig. 13. Scalability evaluation for variable number of HMC devices.

effectively minimizes remote access between multiple HMC cubes while ensuring load balancing. Consequently, the HMC-FHE framework facilitates enhancing the overall parallelism of the architecture by scaling up the number of HMC devices without augmenting the complexity of interactions between HMC devices and GPU devices, or among HMC devices themselves. To quantify this characteristic, we reconfigured the number of HMC devices (e.g., 4, 6, 8, 10) in the architecture to assess the impact of changes in the number of HMCs on the overall acceleration performance of the architecture.

The results shown in Fig. 13 reveal that when there are ample computational resources on both GPU and HMC ends, HMC-FHE exhibits nearly linear performance improvement across all workloads, particularly evident in the cases of HMC\*4 and HMC\*8. For the remaining two cases, the performance enhancement is slightly lower than the proportion of changes in the number of HMC devices. This discrepancy arises mainly because when the number of HMCs is not a power of two, the strides within each stage of the NTT and iNTT kernels cannot consistently remain congruent with the number of HMCs, thus introducing some additional remote accesses between HMCs. In summary, the aforementioned results underscore the robust scalability of HMC-FHE.

### E. Comparison With NVIDIA Tesla V100

To further demonstrate the superiority of HMC-FHE, we used it as a baseline and compared it against the actual GPU device, NVIDIA Tesla V100. We replicated the six workloads mentioned in Section V-A on the V100, with experimental setups consistent with previous tests, and the results are shown in Fig. 14. First, it was observed that across all workloads, the performance of the NVIDIA Tesla V100 consistently lagged behind HMC-FHE. In various workloads, its inference latency achieved only  $0.69\times$ ,  $0.57\times$ ,  $0.69\times$ ,  $0.64\times$ , and  $0.60\times$  that of HMC-FHE. Additionally, the NVIDIA Tesla V100 also demonstrated disadvantages in energy consumption, with  $0.59\times$ ,  $0.67\times$ ,  $0.60\times$ ,  $0.62\times$ , and  $0.65\times$  that of HMC-FHE.

It is noteworthy that, unlike the comparison in Section V-B between the baseline and HMC-FHE where computational power and bandwidth were aligned, the NVIDIA Tesla V100, despite having higher-peak computational power and greater bandwidth, still did not achieve an advantage in inference latency. This further validates our point: due to the different computational and memory access characteristics of various FHE kernels, merely increasing resources does not lead to

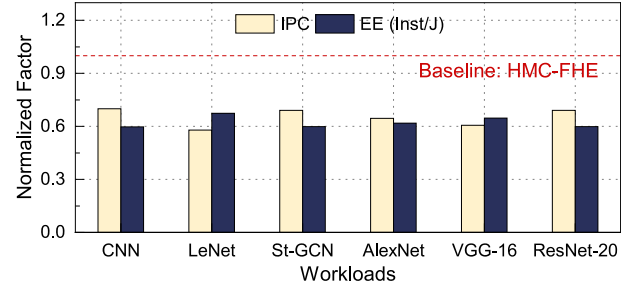


Fig. 14. Performance comparison between the HMC-FHE with four HMC devices and NVIDIA Tesla V100.

a Pareto-optimal performance solution. Appropriate kernel scheduling and resource allocation are far more crucial.

### F. Comparison With Other Platforms

We also conducted comparative experiments to assess the full-system performance of HMC-FHE against other state-of-the-art accelerator prototypes, including FPGA and ASIC. To better reflect the performance of GPUs, we have set the original baseline to the NVIDIA Tesla V100. It is important to note that the computational power and bandwidth resources of the V100 are slightly greater than those of the HMC-FHE (with four HMC devices) architecture. To ensure a fair performance comparison, we standardized the CKKS encryption parameters to  $N = 2^{16}$ ,  $L = 45$  across the aforementioned accelerator prototypes.

1) *Latency*: As shown in Table IV, utilizing an HMC-FHE system with four HMC devices outperforms Poseidon [29], and F1 [5] accelerators in terms of inference latency. However, compared to BTS [7] accelerators equipped with large-capacity SRAM, the inference performance of HMC-FHE is 58% lower. Nonetheless, as previously mentioned, HMC-FHE exhibits superior scalability, and better performance can be achieved by doubling the number of HMC devices.

2) *SRAM Usage*: Concerning SRAM usage, Table IV illustrates that HMC-FHE consumes the least amount of SRAM. Conversely, ASIC-based acceleration designs have reached an impressive usage of 512MB of SRAM, emphasizing their significant demand for SRAM resources. This comparison underscores the advantage of HMC-FHE in resource efficiency, especially in environments with constrained SRAM resources or under heavy workloads (e.g., encryption parameters).



TABLE IV  
COMPARISON OF PERFORMANCE BETWEEN HMC-FHE AND OTHER  
PRIOR WORKS FOR RESNET-20 INFERENCE

	Poseidon (FPGA)	F1+ (ASIC)	BTS (ASIC)	HMC-FHE (HMC*4)	HMC-FHE (HMC*8)
Time (s)	2.67	2.69	1.545	2.44	1.23
Speedup	0.91×	0.91×	1.58×	1×	2×
SRAM (MB)	8.6	256	512	9.25	14.5
Speedup	0.93×	27×	55×	1×	1.56×
Energy (W)	600	93	163	635	640
Speedup	0.95×	0.15×	0.26×	1×	1×

915 3) *Energy Consumption*: For energy consumption, the  
916 results shown in Table IV clearly demonstrate that ASIC  
917 solutions exhibit the lowest. Our proposed approach sig-  
918 nificantly improves upon the power efficiency compared to  
919 GPU deployments (as in Section V-E) and achieves energy  
920 consumption levels similar to those of FPGAs.

## 921 VI. RELATED WORK

922 *ASIC*: ASIC-based FHE accelerators, such as ARK [8],  
923 F1 [5], BTS [7], and CraterLake [16], significantly enhance  
924 the performance of FHE operations. Benefiting from the  
925 flexible resource usage and relatively more freedom in the  
926 design phase, these schemes maximize parallelism and data  
927 movement efficiency during the execution of FHE operations  
928 by customizing specialized processing units and using high-  
929 capacity SRAM.

930 *FPGA*: As a customizing-computing device, the FHE accel-  
931 erator [4], [9], [10], [11], [29], [39] using FPGA adopts  
932 the idea of pipeline to optimize the execution process.  
933 HEAX [10] is a typical FPGA-accelerator in this field,  
934 introducing a pipeline-parallel architecture that maximizes par-  
935 allelism from ciphertext to modular arithmetic logic. Building  
936 on this, FxHENN [4] integrates software/hardware co-design  
937 to achieve efficient resource management.

938 *GPU*: Due to its abundant parallel resources, GPU,  
939 as a general-purpose commercial device, is a promising  
940 scheme to accelerate FHE operations [12], [13], [14], [15].  
941 TensorFHE [12] taps into the computing power of Tensor Core  
942 to accelerate the computational process of CKKS with fine-  
943 grained batch processing measures. From the perspective of  
944 compute-memory balance, 100× [13] avoids the extra memory  
945 access cost by means of operation fusion to speed up the  
946 operations.

947 *PIM*: Some works also adopt *Processing-in-Memory* to  
948 accelerate FHE operations [40], [41]. MemFHE [40] designs  
949 a configurable pipeline to accelerate FHE operations by  
950 mining the features of PIM in-situ computation and extensive  
951 parallelism. Similarly, CryptoPIM [41] optimizes memory-  
952 intensive FHE computing tasks around the excess internal  
953 bandwidth provided by the PIM architecture by building  
954 custom processing units and special memory requirements.  
955 These works focus on deploying FHE operations using only  
956 NDP devices, orthogonal to our work.

## 957 VII. CONCLUSION

958 In this article, we present HMC-FHE, an acceleration frame-  
959 work based on the heterogeneous GPU-HMC architecture

to provide resource management and performance accel- 960  
eration for FHE applications. HMC-FHE aims to achieve 961  
fine-grained optimization of FHE kernels with diverse fea- 962  
tures and offer optimal global task/resource scheduling to 963  
fully exploit the benefits of the GPU-HMC architecture. 964  
Various evaluation results show that compared with the SOTA 965  
GPU-based acceleration schemes, HMC-FHE achieves up to 966  
2.65× performance improvement and reduces 1.81× energy 967  
consumption. 968

## REFERENCES 969

- [1] J. Lin et al., "INSPIRE: In-storage private information retrieval via 970  
protocol and architecture co-design," in *Proc. ISCA*, 2022, pp. 102–115. 971
- [2] R. Ran, N. Xu, W. Wang, G. Quan, J. Yin, and W. Wen, 972  
"CryptoGCN: Fast and scalable homomorphically encrypted graph  
convolutional network inference," in *Proc. Adv. NeurIPS*, vol. 35, 2022, 973  
pp. 37676–37689. 974
- [3] B. Reagen et al., "Cheetah: Optimizing and accelerating homomorphic 975  
encryption for private inference," in *Proc. IEEE Int. Symp. High-  
Perform. Comput. Archit. (HPCA)*, 2021, pp. 26–39. 976
- [4] Y. Zhu, X. Wang, L. Ju, and S. Guo, "FxHENN: FPGA-based accel- 977  
eration framework for homomorphic encrypted CNN inference," in *Proc.  
HPCA*, 2023, pp. 896–907. 978
- [5] N. Samardzic et al., "F1: A fast and programmable accelerator for fully 979  
homomorphic encryption," in *Proc. 54th Annu. IEEE/ACM Int. Symp.  
Microarchit.*, 2021, pp. 238–252. 980
- [6] N. Zhang et al., "NTTU: An area-efficient low-power NTT-uncoupled 981  
architecture for NTT-based multiplication," *IEEE Trans. Comput.*, vol. 982  
69, no. 4, pp. 520–533, Apr. 2020. 983
- [7] S. Kim et al., "BTS: An accelerator for bootstrappable fully homomor- 984  
phic encryption," in *Proc. 49th Annu. Int. Symp. Comput. Archit.*, 2022, 985  
pp. 711–725. 986
- [8] J. Kim et al., "ARK: Fully homomorphic encryption accelerator with 987  
runtime data generation and inter-operation key reuse," in *Proc. MICRO*, 988  
2022, pp. 1237–1254. 989
- [9] Y. Yang, H. Zhang, S. Fan, H. Lu, M. Zhang, and X. Li, "Poseidon: 990  
Practical homomorphic encryption accelerator," in *Proc. IEEE Int. Symp.  
High-Perform. Comput. Archit. (HPCA)*, 2023, pp. 870–881. 991
- [10] M. S. Riazzi, K. Laine, B. Pelton, and W. Dai, "HEAX: An archi- 992  
tecture for computing on encrypted data," in *Proc. ASPLOS*, 2020, 993  
pp. 1295–1309. 994
- [11] S. Sinha Roy, F. Turan, K. Jarvinen, F. Vercauteren, and I. Verbauwhede, 995  
"FPGA-based high-performance parallel architecture for homomorphic  
computing on encrypted data," in *Proc. HPCA*, 2019, pp. 387–398. 996
- [12] S. Fan, Z. Wang, W. Xu, R. Hou, D. Meng, and M. Zhang, "TensorFHE: 997  
Achieving practical computation on encrypted data using GPGPU," 998  
in *Proc. 29th IEEE Int. Symp. High-Perform. Comput. Archit.*, 2023, 999  
pp. 922–934. 1000
- [13] W. Jung, S. Kim, J. H. Ahn, J. H. Cheon, and Y. Lee, "Over 100x faster 1001  
bootstrapping in fully homomorphic encryption through memory-centric  
optimization with GPUs," *IACR Trans. Cryptogr. Hardw. Embed. Syst.*, 1002  
no. 4, pp. 114–148, 2021. 1003
- [14] L. de Castro et al., "Does fully homomorphic encryption need compute 1004  
acceleration?" *Cryptol. ePrint Arch.*, IACR, Bellevue, WA, USA, 1005  
Rep. 2021/1636, 2021. 1006
- [15] W. Jung et al., "Accelerating fully homomorphic encryption through 1007  
architecture-centric analysis and optimization," *IEEE Access*, vol. 9, 1008  
pp. 98772–98789, 2021. 1009
- [16] N. Samardzic et al., "CraterLake: A hardware accelerator for efficient 1010  
computation on encrypted data," in *Proc. 49th Annu. Int. Symp.  
Comput. Archit.*, 2022, pp. 173–187. 1011
- [17] N. Goswami, R. Shankar, M. Joshi, and T. Li, "Exploring GPGPU 1012  
workloads: Characterization methodology, analysis and microarchitec-  
ture evaluation implications," in *Proc. IISWC*, 2010, pp. 1–10. 1013
- [18] H. Jin et al., "Accelerating graph convolutional networks through 1014  
a PIM-accelerated approach," *IEEE Trans. Comput.*, vol. 72, no. 9, 1015  
pp. 2628–2640, Sep. 2023. 1016
- [19] L. Nai, R. Hadidi, J. Sim, H. Kim, P. Kumar, and H. Kim, "GraphPIM: 1017  
Enabling instruction-level PIM offloading in graph computing frame-  
works," in *Proc. IEEE Int. Symp. High Perform. Comput. Archit.  
(HPCA)*, 2017, pp. 457–468. 1018

- 1030 [20] J. Chen, Z. Zhong, K. Sun, C. Ma, R. Mao, and Y. Wang, "Lift: Exploiting hybrid stacked memory for energy-efficient processing of graph convolutional networks," in *Proc. 60th ACM/IEEE Design Autom. Conf. (DAC)*, 2023, pp. 1–6.
- 1031  
1032  
1033
- 1034 [21] J. Chen et al., "GCIM: Toward efficient processing of graph convolutional networks in 3-D-stacked memory," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 41, no. 11, pp. 3579–3590, Nov. 2022.
- 1035  
1036
- 1037 [22] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Proc. Int. Conf. Theory Appl. Cryptol. Inf. Secur.*, 2017, pp. 409–437.
- 1038  
1039
- 1040 [23] J. Fan and F. Vercauteren, "Somewhat practical fully homomorphic encryption," *Cryptol. ePrint Arch., IACR, Bellevue, WA, USA, Rep. 2012/144*, 2012.
- 1041  
1042
- 1043 [24] Z. Brakerski, C. Gentry, and V. Vaikuntanathan, "(Leveled) fully homomorphic encryption without bootstrapping," *ACM Trans. Comput. Theory*, vol. 6, no. 3, pp. 1–36, 2014.
- 1044  
1045
- 1046 [25] K. Hsieh et al., "Transparent offloading and mapping (TOM): Enabling programmer-transparent near-data processing in GPU systems," *ACM SIGARCH Comput. Archit. News*, vol. 44, no. 3, pp. 204–216, 2016.
- 1047  
1048
- 1049 [26] X. Zhang, S. L. Song, C. Xie, J. Wang, W. Zhang, and X. Fu, "Enabling highly efficient capsule networks processing through a PIM-based architecture design," in *Proc. IEEE Int. Symp. High Perform. Comput. Archit. (HPCA)*, 2020, pp. 542–555.
- 1050  
1051  
1052
- 1053 [27] S. Williams, A. Waterman, and D. A. Patterson, "Roofline: An insightful visual performance model for multicore architectures," *Commun. ACM*, vol. 52, no. 4, pp. 65–76, 2009.
- 1054  
1055
- 1056 [28] "NVIDIA Nsight." 2018. [Online]. Available: <https://developer.nvidia.com/nsight-compute>
- 1057
- 1058 [29] Y. Yang, H. Zhang, S. Fan, H. Lu, M. Zhang, and X. Li, "Poseidon: Practical homomorphic encryption accelerator," in *Proc. IEEE Int. Symp. High-Perform. Comput. Archit. (HPCA)*, 2023, pp. 870–881.
- 1059  
1060
- 1061 [30] K. Han and D. Ki, "Better bootstrapping for approximate homomorphic encryption," in *Proc. Cryptogr. Track RSA Conf.*, 2020, pp. 364–390.
- 1062
- [31] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Math. Comput.*, vol. 19, no. 90, pp. 297–301, 1965.
- 1063  
1064  
1065
- [32] M. Khairy, Z. Shen, T. M. Aamodt, and T. G. Rogers, "Accel-sim: An extensible simulation framework for validated GPU modeling," in *Proc. ACM/IEEE 47th Annu. Int. Symp. Comput. Archit. (ISCA)*, 2020, pp. 473–486.
- 1066  
1067  
1068  
1069
- [33] D.-I. Jeon and K.-S. Chung, "CasHMC: A cycle-accurate simulator for hybrid memory cube," *IEEE Comput. Archit. Lett.*, vol. 16, no. 1, pp. 10–13, Jan.–Jun. 2016.
- 1070  
1071  
1072
- [34] Y. Shen, L. Schreuders, A. Pathania, and A. D. Pimentel, "Thermal management for 3-D-stacked systems via unified core-memory power regulation," *ACM Trans. Embed. Comput. Syst.*, vol. 22, no. 5s, pp. 1–26, Sep. 2023.
- 1073  
1074  
1075  
1076
- [35] N. Jiang, G. Michelogiannakis, D. Becker, and B. Towles, *BookSim 2.0 User's Guide*, Stanford Univ., Stanford, CA, USA, 2010, p. q1.
- 1077  
1078
- [36] V. Kandiah et al., "AccelWatch: A power modeling framework for modern GPUs," in *Proc. 54th MICRO*, 2021, pp. 738–753.
- 1079  
1080
- [37] S. H. Pugsley et al., "NDC: Analyzing the impact of 3-D-stacked memory+logic devices on MapReduce workloads," in *Proc. IEEE Int. Symp. Perform. Anal. Syst. Softw. (ISPASS)*, 2014, pp. 190–200.
- 1081  
1082  
1083  
1084
- [38] Y. Cai, Q. Zhang, R. Ning, C. Xin, and H. Wu, "Hunter: HE-friendly structured pruning for efficient privacy-preserving deep learning," in *Proc. AsiaCCS*, 2022, pp. 931–945.
- 1085  
1086  
1087
- [39] M. Jiang et al., "FHE-CGRA: Enable efficient acceleration of fully homomorphic encryption on CGRAs," in presented at DAC, 2024.
- 1088  
1089
- [40] S. Gupta, R. Cammarota, and T. Š. Rosing, "MemFHE: End-to-end computing with fully homomorphic encryption in memory," *ACM Trans. Embed. Comput. Syst.*, vol. 23, no. 2, pp. 1–23, 2024.
- 1090  
1091  
1092
- [41] H. Nejatollahi, S. Gupta, M. Imani, T. S. Rosing, R. Cammarota, and N. Dutt, "CryptoPIM: In-memory acceleration for lattice-based cryptographic hardware," in *Proc. DAC*, 2020, pp. 1–6.
- 1093  
1094  
1095