

Efficient TMVP-based Polynomial Convolution on GPU for Post-quantum Cryptography Targeting IoT Applications

Author's Version

Muhammad Asfand Hafeez*, *Graduate Student Member, IEEE*, Wai-Kong Lee†, *Member, IEEE*, Angshuman Karmakar‡, *Member, IEEE*, and Seong Oun Hwang‡, *Senior Member, IEEE*

Abstract—Recently proposed lattice-based cryptography algorithms can be used to protect the IoT communication against the threat from quantum computers, but they are computationally heavy. In particular, polynomial convolution is one of the most time-consuming operations in lattice-based cryptography. To achieve efficient implementation, the Number Theoretic Transform (NTT) algorithm is an ideal choice, but it has certain limitations on the parameters, which not all lattice-based schemes can employ directly. Hence, alternative techniques are proposed to accelerate polynomial convolution on lattice-based schemes that cannot utilize the NTT directly. In this paper, we propose a parallel Toeplitz matrix-vector product (TMVP) version to accelerate the polynomial convolution in PQC algorithms implemented it on a graphics processing unit (GPU). This is the first time a TMVP parallel version has been proposed and experimented on different GPU cores (i.e., CUDA-cores and Tensor-cores). The effectiveness of the proposed solution is validated on Saber (the NIST post-quantum standardization finalist) and Sable (an improved version of Saber) schemes. Experimental results show that TMVP-based polynomial convolution using CUDA-cores fails to exhibit a significant enhancement compared to the schoolbook CUDA-core method already proposed by Hafeez et al. 2023. However, when the TMVP technique is applied to Tensor-cores, it outperformed state-of-the-art implementations. The proposed Tensor-core approach outperformed the schoolbook Tensor-core method by up to 1.21×, and outperformed the dot-product-instructions method (Lee et al. 2022) by up to 3.63×. The proposed TMVP Tensor-cores is also faster than the TMVP CUDA-cores method by 13.76×.

Index Terms—Toeplitz Matrix-vector Product (TMVP), Cryptography, Tensor-cores, CUDA-cores, Post-quantum Cryptography, Lattice-based Cryptography, Matrix Multiplication.

I. INTRODUCTION

SECURE communication is essential for protecting sensitive information and preserving privacy; it relies heavily on cryptography algorithms to achieve the desired goals. However, the emergence of quantum computers (QCs) poses a significant threat to the security provided by the classical cryptography schemes relying on the hardness of integer

factorization and discrete logarithms. In response to this threat, the National Institute of Standards and Technology (NIST) [1] started a Post-Quantum Cryptography (PQC) standardization process in 2016. The goal was to identify cryptography algorithms that could resist attacks from classical and quantum computers in the long term. After a comprehensive evaluation process, lattice-based algorithms emerged as the most resilient option for PQC. The standardization process concluded in 2022 with four candidates: one key encapsulation mechanism (KEM), CRYSTALS-KYBER [2] and three signature schemes CRYSTALS-Dilithium [3], FALCON [4], and SPHINCS+ [5].

A. Continuous Development of PQC

Although using the Kyber algorithm as the primary standard for PQC is a significant step forward, it provides a framework for future advancements and improvements in PQC schemes while also reinforcing the importance of security. During the standardization process, non-traditional parameter choices, such as the LAC and Round 5 [6], were discouraged to mitigate the potential vulnerabilities that attackers could exploit. This cautious approach ensures that the security of PQC systems remains robust.

The use of non-constant-time error correction codes in lattice-based PQC schemes has raised concerns. Error correction codes play a crucial role in ensuring the accuracy and dependability of PQC schemes. However, if these codes are not implemented in a constant-time manner, they can become potential sources of side-channel attacks. These attacks can compromise the security of the system by exploiting information leaked through timing or power consumption. Therefore, it is essential to evaluate the use of error correction codes in PQC schemes with care. Researchers and developers must continually strive to enhance existing PQC schemes while maintaining their security. This ongoing effort includes exploring alternative parameter choices, optimizing error correction codes, and addressing potential side-channel vulnerabilities, etc. For example, Scabbard (a suite of KEM schemes proposed by Mera et al. [7]) improves on Saber [8], the NIST PQC finalist. SMAUG which is a candidate scheme submitted to the ongoing Korean PQC standardization [9] has been heavily influenced by the design elements of Scabbard. Similarly, Liang et al. [10] proposed an enhanced version of the NTRU KEM [11], which was also a finalist in the NIST

The author* is with the Department of IT Convergence Engineering, Gachon University, Seongnam 13120, South Korea (e-mail: muhammadasfandh@gmail.com) and authors† are with the Department of Computer Engineering, Gachon University, Seongnam 13120, South Korea (e-mail: waikong.lee@gmail.com; sohwang@gachon.ac.kr). The author‡ is with the Department of Computer Science and Engineering, Indian Institute of Technology, Kanpur 208016, India (e-mail: angshuman@cse.iitk.ac.in) and with imec-COSIC, Department Electrotechniek, KU Leuven Belgium (e-mail: angshuman.karmakar@esat.kuleuven.be). (Corresponding author: Seong Oun Hwang.)

standardization. Cho et al. [12] improved the key size and bit-security of the first-round pqsigRM signature scheme.

However, most lattice-based PQC schemes involve polynomial convolution with high-degree polynomials, which makes them computationally expensive. To overcome this, some schemes like Kyber [2] are designed to have a special ring structure that can utilize the Number Theoretic Transform (NTT) [13] for computing polynomial convolution. However, some lattice-based schemes like Saber have power-of-two numbers ($q = 2^{13}$ and $p = 2^{10}$), which cannot natively support NTT. For NTT to be compatible, the inverse of the length of the polynomial, i.e., 256^{-1} , must exist in the field of Z_q or Z_p . Since this is not the case for Saber, we have to optimize polynomial multiplication using different techniques.

On the other hand, Kyber [2] uses a prime number q , which makes it compatible with NTT. NTT is faster than other multiplication techniques, such as Toom-Cook, Karatsuba, or TMVP-based multiplication. However, schemes that use power-of-two fields cannot use NTT, so they have to optimize polynomial multiplication using different techniques. One might argue that we can use TMVP-based multiplication for Kyber-like schemes. This is theoretically fine, but we need to make many modular reductions while using TMVP-based multiplication. This is trivial for power-of-two moduli, but incurs significant overhead for prime modulus.

The efficiency of polynomial convolution/multiplication operations significantly impacts the overall performance and security of the PQC algorithms. Therefore, careful consideration should be given while implementing lattice-based schemes to ensure optimal performance. Substantial efforts have been directed toward enhancing the performance of polynomial convolution for non-NTT-friendly schemes. For instance, classical techniques like Toom-Cook [14] and Karatsuba [15] are commonly used to achieve this. Recently, the Toeplitz matrix-vector product (TMVP) emerged as an alternative method, and its effectiveness was demonstrated in recent works [16], [17]. These studies showed that the TMVP yields promising results in terms of performance and efficiency when compared to the Toom-Cook and Karatsuba methods. However, prior work was only focused on serial versions of the TMVP; it is still unclear if such an approach can perform equally well on a parallel architecture like the graphics processing unit (GPU). This motivated us to investigate the effectiveness of a parallel TMVP to speed up polynomial convolution further.

B. Deployment of PQC for IoT Applications

The deployment of PQC algorithms has become an urgent need due to the threat posed by quantum computing and the "Harvest First, Decrypt Later" strategy [18]. However, this poses a significant challenge when it comes to the communication between sensor nodes, gateway devices, and cloud servers. Each node requires the deployment of the PQC algorithm, and the cloud server and gateway devices must handle large amounts of key encapsulation/decapsulation. One effective way to mitigate this challenge is to use accelerators on the server side to achieve high-throughput KEM. This can be achieved by implementing PQC algorithms on GPUs,

which enhances performance and makes them viable for high-throughput operations required in IoT ecosystems.

For instance, Gupta et al. [19] presented the possibility of using GPU to implement PQC algorithms with high performance. Lee et al. [20] proposed the concept of key encapsulation/decapsulation as a service (KEDaaS), which is very useful to IoT applications. Moreover, the introduction of Tensor-cores in NVIDIA GPUs has further enhanced the potential of these devices in cryptographic contexts. These cores are adept at handling polynomial convolutions, a critical operation in many lattice-based PQC algorithms, thus enabling higher throughput in KEM and KEX. This, in turn, supports IoT applications in smart grids [21], healthcare [22], and industrial IoT [23], where we need to process a multitude of data in a single communication cycle. On the other hand, gateway devices can also employ embedded platforms like Jetson [24] and ODROID [25] that come with GPU to compute the PQC algorithms. Note that GPU is already a de-facto accelerator in many cloud service providers (e.g., AWS) targeting IoT applications, which makes it more favourable compared to other accelerators like ASIC and FPGA.

C. Contributions

In this paper, our primary aim is to investigate the feasibility of parallelizing TMVP to analyze its performance on the GPU platform. We also explore the possibility of utilizing Tensor-cores in conjunction with the TMVP to further improve the performance of polynomial convolution.

- 1) For the first time, TMVP-based polynomial convolution on Tensor-cores on a GPU is presented. This parallel implementation of TMVP presents certain challenges, including memory access patterns, shared memory limitations, and the choice of parallelization methods in order to optimally leverage the capability of GPU architecture. To meet these challenges, we pre-arrange the matrix following the reduction pattern of the selected schemes (Saber and Sable), and then apply the TMVP to break the matrix in a manner that maximizes parallelism. Experimental results on a RTX 3060Ti GPU demonstrate that our proposed TMVP-based polynomial convolution using Tensor-cores yields throughput that is $1.21 \times$ and $3.63 \times$ higher than the [26] and [27], respectively.
- 2) In addition to Tensor-cores, the proposed TMVP-based polynomial convolution was also implemented on CUDA-cores. The findings reveal that the TMVP using Tensor-cores outperformed its CUDA-cores counterpart by $6.2 \times$ in terms of throughput. This is because there is insufficient shared memory to hold multiple copies of vectors in the CUDA-cores TMVP implementation. In addition, many read/write operations are required in the CUDA-cores TMVP implementation, limiting its performance. This shows that the TMVP technique may not always yield good performance in a parallel architecture due to limitations in memory. In contrast, the Tensor-cores version does not use any shared memory because matrix multiplication is performed directly on the registers, thus eliminating most of the memory issues found in the CUDA-cores version.

- 3) The Saber [8] and Sable [7] KEMs were evaluated using the proposed techniques. Our Tensor-cores implementation achieved 424,437 encryptions per second and 6,259,781 decryptions per second implementing the Saber key exchange (KEX) on an RTX 3060Ti GPU, which is $2.58\times$ and $6.83\times$ faster, respectively, than using standard CUDA-cores. The highest throughput achieved by Saber KEM was 267,720 encapsulations per second and 294,020 decapsulations per second. For the Sable KEX, the throughput achieved by the TMVP-based Tensor-cores implementation was 457,155 encryptions per second and 5,621,925 decryptions per second, which is $2.67\times$ and $6.22\times$ faster, respectively, than on standard CUDA-cores. The highest throughput of the Sable KEM was 250,062 encapsulations per second and 295,061 decapsulations per second. The Tensor-core based TMVP implementation for Sable demonstrated satisfactory performance, wherein the encapsulation and decapsulation throughput were 4.7% and 4.97% faster than [26].
- 4) The source code for the proposed TMVP polynomial convolution is publicly available <https://github.com/Muhammad-Asfand/asfand-tmvp>. We sincerely hope that this will enable researchers to easily replicate our findings. Also, we believe that it can encourage further studies and research on TMVP-based polynomial convolution on GPUs and other parallel accelerators.

This paper is organized as follows. Section II discusses background information for the proposed study and reviews related work in the literature. In Section III, we discuss in detail an implementation of the TMVP using CUDA-cores and Tensor-cores. In Section IV, we discuss our experiment results. Finally, Section V concludes the paper.

II. PRELIMINARIES

In this section, we first discuss the deployment of PQC schemes in IoT applications. Then, an overview of the TMVP and its various variants is presented, followed by its applications to reduce the complexity of polynomial convolution for PQC. Two target PQC schemes that can utilize TMVP for improved performance are presented next. The first scheme (Saber) is one of the finalists in the NIST PQC standardization process, and the second scheme (Sable) is the improved version of Saber.

A. The Toeplitz matrix-vector product technique

The TMVP is a technique used in various cryptographic applications to perform multiplication. It was first introduced by Fan and Hasan [28] for multiplying binary extension fields. Since then, many proposals have been suggested by Hasan et al. [29], [30]. Similarly, in [31] and [32], the TMVP was used for speeding up the residue multiplication modulo in integer modular multiplication. It can also be used to calculate the product of two polynomials modulo a polynomial [33]. The following matrix T is an example of a 5×5 Toeplitz matrix where the elements along a line parallel to the principal diagonal possess a constant value.

$$T = \begin{pmatrix} t_0 & t'_1 & t'_2 & t'_3 & t'_4 \\ t_1 & t_0 & t'_1 & t'_2 & t'_3 \\ t_2 & t_1 & t_0 & t'_1 & t'_2 \\ t_3 & t_2 & t_1 & t_0 & t'_1 \\ t_4 & t_3 & t_2 & t_1 & t_0 \end{pmatrix} \quad (1)$$

To determine an $n \times n$ Toeplitz matrix, only $2n - 1$ elements are needed. This means that calculating the sum of two Toeplitz matrices can be done with just $2n - 1$ entry additions, resulting in another Toeplitz matrix. Additionally, all submatrices of a Toeplitz matrix are also Toeplitz matrices. These characteristics make it possible to efficiently compute Toeplitz matrix-vector multiplication using TMVP formulas rather than the conventional schoolbook method.

1) *TMVP Formulas*: Various split formulas are available to efficiently compute TMVPs (such as two-way, three-way, and four-way, given in [31], [16], [17]). We use X to denote an $n \times n$ Toeplitz matrix and Y to denote a vector of length n .

Two-way TMVP (TMVP-2): We can define $n \times n$ Toeplitz matrix T using three matrix vectors, (X_0, X_1, X_2) and an $n \times 1$ column vector $Y = (Y_0, Y_1)$. Toeplitz matrix T consists of three $(n/2) \times (n/2)$ Toeplitz matrices, namely P_0, P_1 , and P_2 . Equation 2 is the TMVP-2 using three $(n/2) \times (n/2)$ TMVPs [34].

$$T = X \cdot Y = \begin{pmatrix} X_1 & X_0 \\ X_2 & X_1 \end{pmatrix} \begin{pmatrix} Y_0 \\ Y_1 \end{pmatrix} = \begin{pmatrix} P_0 + P_1 \\ P_0 - P_2 \end{pmatrix}, \quad (2)$$

where P_0, P_1 and P_2 represents three TMVPs:

$$\begin{aligned} P_0 &= X_1(Y_0 + Y_1), \\ P_1 &= (X_0 - X_1)Y_1, \\ P_2 &= (X_1 - X_2)Y_0. \end{aligned}$$

Three-way TMVP (TMVP-3): Like TMVP-2, TMVP-3 allows us to calculate an n dimensional TMVP using six $n/3$ -dimensional TMVPs. Consider the $n \times 1$ column vector $Y = (Y_0, Y_1, Y_2)$ and matrix-vector $X = (X_0, X_1, X_2, X_3, X_4)$, which is an $n \times n$ Toeplitz matrix. Here, Y_i (where $i = 0, 1, 2$) is an $(n/3) \times 1$ column vector, and X_i (where $i = 0, 1, 2, 3, 4$) is an $(n/3) \times (n/3)$ Toeplitz matrix [34]. By rewriting product $P = XY$, we get equation 3:

$$X \cdot Y = \begin{pmatrix} X_2 & X_1 & X_0 \\ X_3 & X_2 & X_1 \\ X_4 & X_3 & X_2 \end{pmatrix} \begin{pmatrix} Y_0 \\ Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} P_0 + P_3 + P_4 \\ P_1 - P_3 + P_5 \\ P_2 - P_4 + P_5 \end{pmatrix}, \quad (3)$$

where P_0, P_1, P_2, P_3, P_4 and P_5 represents six TMVPs:

$$\begin{aligned} P_0 &= (X_0 + X_1 + X_2)Y_2, \\ P_1 &= (X_1 + X_2 + X_3)Y_1, \\ P_2 &= (X_2 + X_3 + X_4)Y_0, \\ P_3 &= X_1(Y_1 - Y_2), \\ P_4 &= X_2(Y_0 - Y_2), \\ P_5 &= X_3(Y_0 - Y_1). \end{aligned}$$

Four-way TMVP (TMVP-4): To compute an n -dimensional TMVP, a TMVP-4 formula was proposed in [17]. This utilizes a combination of seven $n/4$ -dimensional TMVPs. Assuming that n is divisible by four, we take the

$n \times 1$ column vector $Y = (Y_0, Y_1, Y_2, Y_3)$ and a matrix-vector $X = (X_0, X_1, X_2, X_3, X_4, X_5, X_6)$, which represents an $n \times n$ Toeplitz matrix. In this case, Y_i (where $i = 0, 1, 2, 4$) is an $n/4 \times 1$ column vector, and X_i (where $i = 0, 1, 2, 3, 4, 5, 6$) is an $n/4 \times n/4$ Toeplitz matrix. We divide the Toeplitz matrix and the vector, then compute the product as in equation 4:

$$x.Y = \begin{pmatrix} X_3 & X_2 & X_1 & X_0 \\ X_4 & X_3 & X_2 & X_1 \\ X_5 & X_4 & X_3 & X_2 \\ X_6 & X_5 & X_4 & X_3 \end{pmatrix} \begin{pmatrix} Y_0 \\ Y_1 \\ Y_2 \\ Y_3 \end{pmatrix} = \begin{pmatrix} P_1 - P_2 + 8P_3 - 8P_4 + 27P_5 + P_6 \\ P_1 + P_2 + 4P_3 + 4P_4 + 9P_5 \\ P_1 - P_2 + 2P_3 - 2P_4 + 3P_5 \\ P_0 + P_1 + P_2 + P_3 + P_4 + P_5 \end{pmatrix}, \quad (4)$$

where $P_0, P_1, P_2, P_3, P_4, P_5$ and P_6 represents six TMVPs:

$$\begin{aligned} P_0 &= \frac{1}{12}(12X_6 + 4X_5 - 15X_4 + 5X_3 + 3X_2 - X_1)Y_0, \\ P_1 &= \frac{1}{12}(12X_5 + 8X_4 - 7X_3 - 2X_2 + X_1)(Y_0 + Y_1 + Y_2 + Y_3), \\ P_2 &= \frac{1}{24}(-12X_5 + 16X_4 - X_3 - 4X_2 + X_1)(Y_0 - Y_1 + Y_2 - Y_3), \\ P_3 &= \frac{1}{24}(-6X_5 - X_4 + 7X_3 + X_2 - X_1)(Y_0 - 2Y_1 + 4Y_2 - 8Y_3), \\ P_4 &= \frac{1}{120}(6X_5 - 5X_4 - 5X_3 + 5X_2 - X_1)(Y_0 - 2Y_1 + 4Y_2 - 8Y_3), \\ P_5 &= \frac{1}{120}(4X_5 - 5X_3 + X_1)(Y_0 + 3Y_1 + 9Y_2 + 27Y_3), \\ P_6 &= (-12X_5 + 4X_4 + 15X_3 - 5X_2 - 3X_1 + A_0)Y_3, \end{aligned}$$

Table I presents the arithmetic complexity of the TMVP-2, TMVP-3, and TMVP-4 formulas. The expressions in the table exhibit a recursive nature and reflect the number of operations required to compute the TMVP for a given size n . These expressions are defined in terms of the operations involved in smaller sizes. It is worth noting that, despite TMVP-4 breaking down the n into a smaller matrix compared to the others, TMVP-2 has the smallest recursive term coefficient (i.e., 3) among them, hence exhibiting the lowest arithmetic complexity of the three methods.

TABLE I
ARITHMETIC COMPLEXITY OF TMVP FORMULAS

TMVP's	Arithmetic Complexity
TMVP-2	$M_{TMVP-2}(n) = 3M(n/2) + 3n - 1$
TMVP-3	$M_{TMVP-3}(n) = 6M(n/3) + 5n - 1$
TMVP-4	$M_{TMVP-4}(n) = 7M(n/4) + 11n - 1$

2) *TMVP vs Toom-Cook*: TMVP and Toom-Cook-based multiplications are specialized techniques for optimizing polynomial convolutions and exhibit notable similarities. Nevertheless, the selection of an appropriate method depends heavily on the distinct computational context and hardware prerequisites. When considering the utilization of GPUs, it is recommended to opt for TMVP due to the following reasons.

Parallelism: TMVP allows for a high level of parallelism, which is a key optimization strategy on GPUs. Different parts of the vector can be multiplied in parallel with various sliding windows of the matrix, increasing function throughput.

Memory Bandwidth: TMVP can lead to reduced memory bandwidth requirements compared to Toom-Cook-based polynomial multiplications. Toom-Cook algorithms involve more complex operations and require more memory transfers, which cause a bottleneck on GPUs, especially for large polynomials.

Data Locality: Toeplitz matrices exhibit a discernible pattern where each diagonal that descends from left to right maintains a constant value. This structure efficiently manages memory coherence when storing and manipulating matrices, particularly in GPUs that support coalesced memory access.

B. Saber and Sable

Saber is a lattice-based KEM that relies on module lattices. Saber stands out for its unique feature of polynomial convolution without the use of NTT, which can be daunting for lattice-based cryptography. This approach of Saber has inspired other cryptography schemes like [7], [35] to adopt similar methods. It was named a finalist in the third round of the NIST PQC standardization competition, indicating its potential as a leading solution in cryptographic security. The strength of Saber's security relies on the conjectural hardness of the Module Learning with Rounding (MLWR) problem [36]. The security level of the target schemes can be configured by specifying dimension ℓ of the module with three distinct values: $\ell = 2$ (LightSaber), $\ell = 3$ (Saber), and $\ell = 4$ (FireSaber), which correspond to security levels 1, level 3, and level 5, respectively. Note that in this paper, we focus on our implementation of Saber for $\ell=3$, extending it to support different ℓ levels is straightforward. Saber's arithmetic operations are $R_q = R_{2^{13}} = \mathbb{Z}_{2^{13}}[x]/\langle x^{256} + 1 \rangle$ and $R_p = R_{2^{10}} = \mathbb{Z}_{2^{10}}[x]/\langle x^{256} + 1 \rangle$. As with many lattice-based cryptosystems defined on polynomial rings, the efficiency of this scheme is heavily impacted by multiplication in these rings. However, it is important to note that the rings R_q and R_p utilized by Saber are not directly compatible with the NTT, which is currently the most efficient polynomial convolution algorithm known.

Mera et al. [7] introduced the Sable scheme in Scabbard as an improved version of Saber based on a hard lattice problem known as learning with rounding (LWR). In such schemes, errors are implicitly created through rounding instead of explicit addition, as seen in LWE. Since errors are crucial in determining the security of lattice-based schemes, proper estimation is essential to avoid overestimation or underestimation. By accurately estimating errors, Mera et al. [7] were able to enhance Saber's parameters without compromising its security. This resulted in reduced key sizes and bandwidth, and this improved version of Saber is known as Sable. The security level of Sable can be configured in the same way as Saber. For instance, $\ell = 2$ (LightSable), $\ell = 3$ (Sable), and $\ell = 4$ (FireSable), correspond to security levels 1, 3, and 5, respectively. The Saber and Sable KEMs consist of three algorithms: key generation (Algorithm 1), encapsulation (Algorithm 2), and decapsulation (Algorithm 3). The values of different parameters used in the designing of both KEMs are given in Table II.

Algorithm 1 depicts the generation of a public key (PK) and a private key (SK) using security parameter N . Algorithm 2 takes the PK as input and produces ciphertext (CT) and a shared secret key (K). Algorithm 3 performs decapsulation, taking the PK, CT, and SK as input and returning the shared secret key as output. In Algorithms 1 to 3, H and \mathcal{G} represent hash functions. The constant polynomials $h1, h2$, and $h3$ have coefficients of $2^{(\epsilon_q - \epsilon_p - 1)}$, $(2^{(\epsilon_q - \epsilon_p - 1)} + 2^{(\epsilon_q - B - 1)} - 2^{(\epsilon_q - \epsilon_t - 1)})$ and $2^{(\epsilon_q - \epsilon_p - 1)}$, respectively.

C. Related work

Recent advancements showcased significant efforts towards optimizing the performance of PQC algorithms through vari-

Algorithm 1 KEM Key Genreation**Data:** nil**Result** PK = ($seed_A$, b), SK = (s, $H(\text{PK})$, r)

- 1: $seed_A \leftarrow \mathcal{U}(\{0, 1\}^{256})$
- 2: $r \leftarrow \mathcal{U}(0, 1)^{256}$
- 3: $A \leftarrow gen_N^{L \times L}(\text{XOF}(seed_A)) \in (\mathcal{R}_q^N)^{L \times L}$
- 4: $s \leftarrow \beta_n((\mathcal{R}_q^N)^L)$
- 5: $b = \text{bits}(A \cdot s + h_1, \epsilon_q, \epsilon_p) \in (\mathcal{R}_q^N)^L$
// Rounding
- 6: $\text{PK} \leftarrow (seed_A, b)r \leftarrow_{\$} \{0, 1\}^{256}$
- 7: $\text{SK} \leftarrow (s, H(\text{PK}), r)$
- 8: **return**
- 9: $\text{PK} = (seed_A, b)$, $\text{SK} = (s, H(\text{PK}), r)$

Algorithm 2 KEM Encapsulation**Data:** PK = ($seed_A$, b)**Result** CT = (c' , b'), key = K

- 1: $m' \leftarrow_{\$} \{0, 1\}^{256}$
- 2: $m = \text{arrange_msg}(m')$
- 3: $(K', r') \leftarrow \mathcal{G}(m || H(\text{PK}))$
- 4: $r' \leftarrow \mathcal{U}(\{0, 1\}^{256})$
- 5: $A \leftarrow gen_N^{L \times L}(\text{XOF}(seed_A)) \in (\mathcal{R}_q^N)^{L \times L}$
- 6: $s' \leftarrow \beta_\eta((\mathcal{R}_q^N)^L)$
- 7: $b' = \text{bits}(A^T \cdot s' + h_1, \epsilon_q, \epsilon_p)$
// Rounding
- 8: $u' = b^T \cdot (s' \bmod p) \in \mathcal{R}_p^N$
- 9: $c' = \text{bits}((u' + h_3 - 2^{\epsilon_p - B} m), \epsilon_p, (\epsilon_t + B)) \in \mathcal{R}_{2^{B_t}}^N$ \triangleright
HelpDecode
- 10: $K \leftarrow H(K', H(c'))$
- 11: **return**
CT = (c' , b'), key = K

Algorithm 3 KEM Decapsulation**Data:** PK = ($seed_A$, b), SK = (s, $H(\text{PK})$, r), CT = (c' , b')**Result** key = K

- 1: $u = b' \cdot (s \bmod p) \in \mathcal{R}_p^N$
- 2: $m'_1 = \text{bits}((u + h_2 - 2^{\epsilon_p - \epsilon_t - B} m), \epsilon_p, B) \in \mathcal{R}_{2^B}^N$ \triangleright
Decode
- 3: $m_1 = \text{original_msg}(m'_1)$
- 4: $m_2 = \text{arrange_msg}(m_1)$
- 5: $(K'_1, r'_1) \leftarrow \mathcal{G}(m_2 || H(\text{pk}))$
- 6: $A \leftarrow gen_N^{L \times L}(\text{XOF}(seed_A)) \in (\mathcal{R}_q^N)^{L \times L}$
- 7: $s'_1 \leftarrow \beta_\eta((\mathcal{R}_q^N)^L)$
- 8: $b'_1 = \text{bits}(A^T \cdot s'_1 + h_1, \epsilon_q, \epsilon_p)$
// Rounding
- 9: $u'_1 = b'^T \cdot (s'_1 \bmod p) \in \mathcal{R}_p^N$
- 10: $c'_1 = \text{bits}((u'_1 + h_3 - 2^{\epsilon_p - B} m), \epsilon_p, (\epsilon_t + B)) \in \mathcal{R}_{2^{B_t}}^N$ \triangleright
HelpDecode
- 11: **if** $c' = c'_1$ **then**
- 12: **return** $K = H(K'_1, H(c'))$
- 13: **else**
- 14: **return** $K = H(r, H(c'))$
- 15: **end if**

a polynomial restructuring technique that enables multiple polynomials with different public keys to be processed in a single communication cycle. Secondly, they introduce a new method to handle the reduction patterns that are not suitable for parallel implementation.

Gao et al. [39], and Lee and Hwang [20] delved into the application of the Number Theoretic Transform (NTT) on GPUs. They implement NTT on a GPU for NewHope and Kyber PQC algorithms, respectively. These studies revealed the potential for GPUs to effectively handle high throughput PQC computation, which is crucial for securing the IoT communication on the server side.

In addition to the above, researchers have also investigated the GPU based implementations of other cryptography schemes. Sun et al. [40] demonstrated an efficient parallel implementation of SPHINCS on a GPU, while Dai et al. [41] optimized the NTRU modular lattice signature scheme for parallel polynomial convolution on a GPU. Their optimization is particularly important due to the scheme's reliance on large vectors, which can be efficiently processed in parallel on a GPU. Additionally, Gupta et al. [19] analyzed the batch mode and single mode parallelism available in a GPU and evaluated implementation in different PQC schemes. The findings of these studies shed light on the potential of utilizing a GPU to provide efficient and scalable solutions for various cryptographic applications.

Similarly, efforts have been made to improve the performance of lattice-based schemes on several alternative platforms such as the latest Intel AVX [42] instructions, hardware accelerators in a Field Programmable Gate Array (FPGA) [43], [44], reduced instruction set computer (RISC) [45] or an application-specific integrated circuit (ASIC) [46] platform. These prior works are applicable to IoT applications to protect the device-side security.

ous computational approaches. In recent work, Lee et al. [27] introduced a novel approach to conduct polynomial convolution using dot-product instructions, which enables the simultaneous execution of *MULTIPLY*-and-*ADD* operations within a single clock cycle. This innovative technique marks a considerable enhancement in throughput compared to the traditional methodologies relying on 32-bit integer units. In addition to this, Lee et al. [37] further explored the utilization of Tensor-cores within GPUs to compute polynomial convolution, which showed greater efficiency and speed compared to CUDA-cores. Following this, Hafeez et al. [26] introduced two techniques to address gaps in previous research. First, they extended the work of See et al. [38] on a GPU and proposed

TABLE II
PARAMETERS OF SABER AND SABLE

Parameters	ℓ	N	p	q	Moduli	Key Sizes
Saber	3	256	2048	8192	$\epsilon_q:13$	PK: 992
					$\epsilon_p:10$	SK: 1440
					$\epsilon_t:4$	CT: 1088
Sable	3	256	512	2048	$\epsilon_q:11$	PK: 1280
					$\epsilon_p:9$	SK: 1728
					$\epsilon_t:4$	CT: 1304

III. PROPOSED PARALLEL TMVP TECHNIQUE

In this section, we describe how to parallelize the TMVP-2 formula and its implementation for Saber and Sable, using Tensor-cores and CUDA-cores.

A. Polynomial convolution using TMVP-2

Saber and Sable schemes both employ an efficient reduction pattern that resembles a nega-cyclic convolution. To facilitate matrix-vector multiplication, the matrix dimension is set to 256×256 , as specified in Table II, the polynomial degree N is 256. This is achieved by transforming the polynomial A into the nega-cyclic matrix as given in equation 5, which yields a matrix of dimension 256×256 . Polynomial B is structured into the column-major matrix in equation 6.

$$A = \begin{pmatrix} a_0 & -a_{n-1} & -a_{n-2} & \dots & -a_3 & -a_2 & -a_1 \\ a_1 & a_0 & -a_{n-1} & \dots & -a_4 & -a_3 & -a_2 \\ a_2 & a_1 & a_0 & \dots & -a_5 & -a_4 & -a_3 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ a_{n-3} & a_{n-4} & a_{n-4} & \dots & a_0 & -a_{n-1} & -a_{n-2} \\ a_{n-2} & a_{n-3} & a_{n-4} & \dots & a_1 & a_0 & -a_{n-1} \\ a_{n-1} & a_{n-2} & a_{n-3} & \dots & a_2 & a_1 & a_0 \end{pmatrix} \quad (5)$$

$$B = \begin{pmatrix} b_0 \\ b_1 \\ b_2 \\ \vdots \\ b_{n-3} \\ b_{n-2} \\ b_{n-1} \end{pmatrix} \quad (6)$$

Figures 1, 2, and 3 show polynomial convolution using TMVP-2, TMVP-3, and TMVP-4 respectively. We opted for TMVP-2 to accelerate the polynomial convolution in Saber and Sable. The rationality analysis for opting for TMVP-2 for polynomial convolution is as follows.

- 1) TMVP-2 break the 256×256 matrix into three non-identical 128×128 matrices as shown in Figure 1. Similarly, TMVP-3 and TMVP-4, respectively, produce five and seven non-identical matrices at 86×86 and 64×64 as depicted in Figure 2 and 3. The matrix size of TMVP-2 is bigger than the other two, so it provides more parallelism than TMVP-3 and TMVP-4.
- 2) Additionally, TMVP-2 performs only the three multiplication in equation 2 to compute the polynomial convolution (see Figure 1), while TMVP-3 and TMVP-4 required six (equation 3) and seven (equation 4) multiplications, respectively.
- 3) TMVP-3 is unsuitable for use in Saber and Sable because the polynomial convolution in these schemes has a length of 256, which is not divisible by 3. Hence, we have to create a 258×258 matrix (divisible by 3) and pad the unused rows and columns with zeroes. After padding, we can perform polynomial convolution using TMVP-3, but there will be some unused rows and columns that waste computational bandwidth.

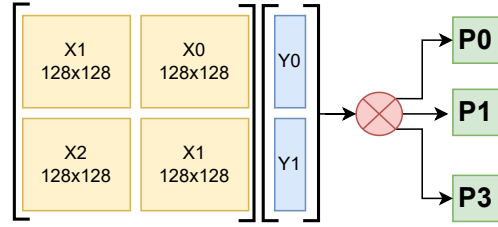


Fig. 1. Polynomial convolution using TMVP-2

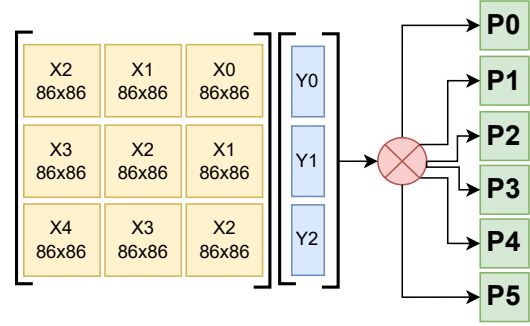


Fig. 2. Polynomial convolution using TMVP-3

B. TMVP-2 Implementation using CUDA-cores

For most of the lattice-based cryptography schemes, polynomial convolution is the most time-consuming task. This particular task entails the manipulation of two distinct polynomials: polynomial a , which typically represents a public or a private key and polynomial b consisting of random elements with small coefficients. In the Sable cryptography algorithms, polynomial b is ternary, i.e., composed of elements $b = \{-1, 0, 1\}$.

However, polynomial convolution in Saber and Sable is essentially the same, so we present the proposed TMVP implementation for both schemes in Algorithm 4. Note that this algorithm describes the basic implementation of the TMVP for polynomial convolution using CUDA-cores commonly found in a GPU. In the next subsection, we present the more advanced technique proposed in this work, which utilizes the Tensor-cores. Referring to Algorithm 4, line 1 rearranges polynomial A into a nega-cyclic pattern. Following this, line 2 pre-processes polynomials A and B for the three TMVP multiplications, as given in equation 2. Next, line 3 computes the matrix-vector product using CUDA-cores, as given in Algorithm 7. Finally, line 7 post-processes the products and calculates the final result.

Algorithm 5 is used to convert the polynomial A into a nega-cyclic pattern. The input is read by N threads and N blocks. Line 3 yields the difference between threads and blocks to arrange the elements into a nega-cyclic pattern. In line 5, if $(tid - bid)$ is greater than the $(N-1)$, the arranged elements in the rows are converted to negative form. Otherwise, the elements are arranged without conversion.

In reference to Algorithm 6, it pre-processes the polynomial A and B into the required matrices and vectors to perform three TMVP multiplications in CUDA-cores. $N/2$ threads and $N/2$ blocks are launched in parallel. Lines 4 and 5 rearrange the

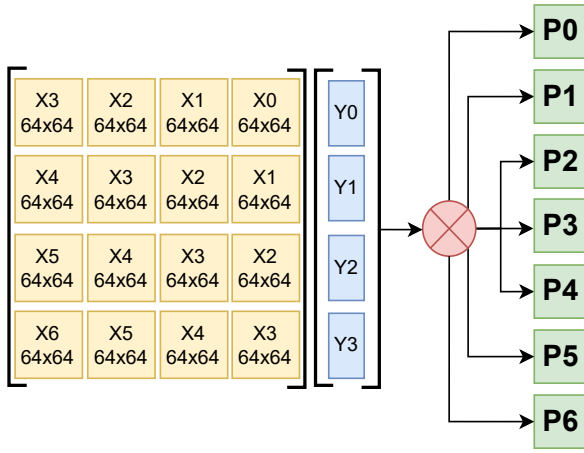


Fig. 3. Polynomial convolution using TMVP-4

Algorithm 4 CUDA-cores implementation of polynomial convolution in parallel on a GPU

Input: Polynomial A , polynomial B , modulus p

Output: $2M \times M$ Matrix c holds the nega-cyclic convolution of polynomial a with polynomial b .

- 1: ParNegCyc $\langle N, N \rangle$ ($fp16_A, A$) ▷ Alg.5
- 2: PreArr $\langle N/2, N/2 \rangle$ ($fp16_B, B$) ▷ Alg.6
- 3: CUDACores $\langle N, N \rangle$ ($fp16_A, fp16_B, fp32_C$) ▷ Alg.7
- 4: PostProcess $\langle N/2, N/2 \rangle$ ($c, fp32_C$) ▷ Alg.8

elements for the first multiplication and move the elements into a_1 and b_1 in U16 format. Similarly, lines 6 and 7 rearrange the elements for the second multiplication, and then lines 8 and 9 rearrange the elements for the third and store the output in a_2, b_2 , and in a_3, b_3 in U16 format, respectively.

After pre-processing, Algorithm 7 describes the proposed method to execute the three TMVP multiplications. This is a crucial step in achieving accurate and efficient results in matrix-vector products. To compute the matrix-vector product, N threads are launched in parallel, ensuring more parallelism is exploited. It is worth noting that the three input matrices used in this process are denoted a_1, a_2 , and a_3 , while the vectors are denoted b_1, b_2 , and b_3 . In lines 3-6 of Algorithm 7, elements of a_1 and b_1 are loaded into the shared memory.

Loading elements into shared memory is crucial for efficient implementation. Line 8 initializes the register to accumulate the product, while lines 9-11 compute the matrix-vector product. Finally, at line 12, the result is moved from the register to p_1 , indicating that one TMVP multiplication is completed. The second and third multiplications are done in lines 14-23 and 25-34, respectively, following a similar approach.

It is worth noting that we can compute the matrix-vector product by launching different numbers of threads (e.g., 128, 256, 512, and 1024) to increase the parallelism. However, there is very little increase in throughput because we cannot make multiple copies of vectors due to the limited amount of shared memory. The combined L1 cache and shared memory in the RTX 3060Ti GPU is 128KB (131,072 bytes). From

Algorithm 5 ParNegCyc: arrange polynomial A into a nega-cyclic pattern

Input: N -length polynomial in

Output: Matrix out of $N \times N$ dimensions, with a polynomial arranged in a nega-cyclic pattern.

- 1: $tid = \text{thread ID}$
- 2: $bid = \text{block ID}$
- 3: $idx = tid - bid$
// Launch N blocks and N threads in
// parallel
- 4: **if** $tid < N$ **then**
- 5: **if** $idx > (N - 1)$ **then**
- 6: $out[bid + tid \times N] = in[(idx) \% N] \times (-1)$
- 7: **else**
- 8: $out[bid + tid \times N] = in[(idx) \% N]$
- 9: **end if**
- 10: **else**
- 11: $out[bid + tid \times N] = 0$
- 12: **end if**

Algorithm 6 PreArr: Pre-arrangements of elements for matrix-vector product.

Input: $N \times N$ -length polynomial in_1 and N -length polynomial in_2

Output: Matrix a_1, a_2, a_3 and vector b_1, b_2, b_3 in U16 format

- 1: $tid = \text{thread ID}$
- 2: $bid = \text{block ID}$
// Launch $N/2$ blocks and $N/2$ threads in
// parallel
- 3: **if** $tid < N$ **then**
- 4: $a_1[bid \times N/2 + tid] = in_1[bid \times N/2 + tid]$
- 5: $b_1[tid] = in_2[tid] + in_2[N/2 + tid]$
- 6: $a_2[bid \times N/2 + tid] = in_1[bid \times N/2 + (N \times N/4) + tid] - in_1[bid \times N/2 + tid]$
- 7: $b_2[tid] = in_2[N/2 + tid]$
- 8: $a_3[bid \times N/2 + tid] = in_1[bid \times N/2 + tid] + in_1[bid \times N/2 + (N \times N/2) + tid]$
- 9: $b_3[tid] = in_2[tid]$
- 10: **else**
- 11: $a_1, a_2, a_3[bid \times N/2 + tid] = 0$
- 12: $b_1, b_2, b_3[tid] = 0$
- 13: **end if**

this, shared memory allocations remain limited to 48KB (49,152 bytes) to maintain architectural compatibility. But each element in the matrix and vector is represented using 16 bits (two bytes). The number of elements for both matrix and vector is $128 \times 128 = 16,384$. The total number of elements combined in both matrix and vector is $16,384 \times 2 = 32,768$. So, the total memory required by both matrix and vector is $32,768 \times 2 = 65,536$ bytes (64KB). If we increase the shared memory to 96 KB, the L1 cache will be reduced. This reduction could increase memory access latency for operations not utilizing shared memory, as less cache would be available to store frequently accessed data. This can result in decreased performance for tasks that rely heavily on L1 cache efficiency.

Algorithm 7 Polynomial convolution of the $n/2$ -dimensional TMVP using 256 threads

Input: $N/2 \times N/2$ -length Matrix a_1, a_2, a_3 and N -length vector b_1, b_2 and b_3

Output: N -length vectors, p_1, p_2 , and p_3

```

1:  $tid = \text{thread ID}$ 
2:  $bid = \text{block ID}$ 
   // Copy elements into shared memory for
   // 1st TMVP  $p_1$  in parallel
3: for  $k$  from 0 to  $N/4$  do
4:    $a\_shared[tid + k \times (N)] = a_1[tid + k \times (N)]$ 
5: end for
6:  $b\_shared[tid \times N] = b_1[tid \times N]$ 
7:  $\_\_synctreads()$   $\triangleright$  Synchronize all the threads
   // Accumulate each column in parallel
   // with  $N$  threads
8:  $sum1 = 0$   $\triangleright$  Use register to accumulate
9: for  $i$  from 0 to  $N/4$  do
10:   $sum1 += a\_shared[tid \times (N/4) + i] \times$ 
     $b\_shared[(tid \% 2) \times (N/4) + i]$ 
11: end for
12:  $p_1[bid + tid] = sum1$ 
13:  $\_\_synctreads()$   $\triangleright$  Synchronize all the threads
   // Copy elements into shared memory for
   // 2nd TMVP  $p_2$  in parallel
14: for  $k$  from 0 to  $N/4$  do
15:   $a\_shared[tid + k \times (N)] = a_2[tid + k \times (N)]$ 
16: end for
17:  $b\_shared[tid \times N] = b_2[tid \times N]$ 
18:  $\_\_synctreads()$ 
19:  $sum2 = 0$ 
20: for  $i$  from 0 to  $N/4$  do
21:   $sum2 += a\_shared[tid \times (N/4) + i] \times$ 
     $b\_shared[(tid \% 2) \times (N/4) + i]$ 
22: end for
23:  $p_2[bid + tid] = sum2$ 
24:  $\_\_synctreads()$   $\triangleright$  Synchronize all the threads
   // Copy elements into shared memory for
   // 3rd TMVP  $p_3$  in parallel
25: for  $k$  from 0 to  $N/4$  do
26:   $a\_shared[tid + k \times (N)] = a_3[tid + k \times (N)]$ 
27: end for
28:  $b\_shared[tid \times N] = b_3[tid \times N]$ 
29:  $\_\_synctreads()$ 
30:  $sum3 = 0$ 
31: for  $i$  from 0 to  $N/4$  do
32:   $sum3 += a\_shared[tid \times (N/4) + i] \times$ 
     $b\_shared[(tid \% 2) \times (N/4) + i]$ 
33: end for
34:  $p_3[bid + tid] = sum3$ 

```

Algorithm 8 PostProcess: Parallel algorithm to process the polynomial coefficients via three N -dimensional TMVPs and modulo p

Input: N -length vectors, p_1, p_2 and p_3

Output: Matrix out of N -length degree, with elements in U16 format and modulo p

```

1:  $tid = \text{thread ID}$ 
2:  $bid = \text{block ID}$ 
   // Launching  $N$  threads at maximum
3: if  $tid < N$  then
4:    $out[bid + tid] += (p_1[bid + (tid \times 2)] + p_1[bid + (tid \times$ 
     $2) + 1] + p_2[bid + (tid \times 2)] + p_2[bid + (tid \times 2) + 1]) \% p$ 
5:    $out[bid + tid] += (p_1[bid + (tid \times 2)] + p_1[bid + (tid \times$ 
     $2) + 1] - p_3[bid + (tid \times 2)] - p_3[bid + (tid \times 2) + 1]) \% p$ 
6: else
7:    $out[bid + tid] = 0$ 
8: end if

```

A reduced L1 cache may also lead to increased cache misses, forcing the GPU to access slower, higher-level memory more often, further impacting overall performance. Balancing the L1 cache and shared memory is crucial to optimize both memory-intensive and compute-intensive tasks on the GPU. Therefore, lines 10, 21, and 32 perform the modulus operation to find the exact element on the vector side. The modulus values for 128, 256, 512, and 1024 are 1, 2, 4, and 8, respectively. Nevertheless, it is imperative to note that this operation hinders the multiplication process and ultimately decreases throughput.

Moreover, the shared memory restriction mandates that we can only load into shared memory the elements of one multiplication at a time, preventing us from performing three multiplications in parallel. This limits to performing one multiplication at a time. As a result, this factor adds another drawback to the implementation of TMVP in CUDA-cores, which significantly decreases throughput. However, if a GPU increases shared memory in the future, it could be possible to overcome these limitations. With more shared memory, we could potentially have access to more data, which helps to store multiple copies of vectors and would allow us to perform three multiplication in parallel. This, in turn, would lead to faster processing and improved overall performance.

After computing the matrix-vector product, we need to do some post-processing to get the final result. In Algorithm 8, input polynomials are read by $N/2$ threads in parallel. Lines 4 and 5 show the post-processing steps given in equation 2 and we then perform the modulo p to get the final result.

C. TMVP-2 implementation using Tensor-cores

The implementation of TMVP on CUDA-cores can be improved by utilizing Tensor-cores, the technique for which is presented in Algorithm 9. Lines 1, 2, and 3 in Algorithm 9 calculate the required numbers of threads and blocks to perform multiplication in Tensor-cores. Line 4 rearranges polynomial A into the same nega-cyclic pattern discussed in Section III-A and described in Algorithm 5. After this, line

Algorithm 9 Tensor-cores implementation of polynomial convolution in parallel on the GPU

Input: Polynomial A , polynomial B , modulus $p||q$

Output: $2M \times M$ Matrix c holds the nega-cyclic convolution of polynomial a with polynomial b .

```

// Calculate total number of threads
// required
1:  $threads\_tot = 32 \times 2 \times (N/32)^2$ 
// Calc. number of blocks
2:  $tc\_blocks = threads\_tot / max\_threads$ 
// Number of thread
3:  $tc\_threads = max\_threads$ 
4: ParNegCyc $\langle N, N \rangle (fp16\_A, A)$   $\triangleright$  Alg.5
5: ParU16toFP16 $\langle N/2, N/2 \rangle (fp16\_B, B)$   $\triangleright$  Alg.10
6: TensorCore $\langle tc\_blocks, tc\_threads \rangle$ 
   ( $fp16\_A, fp16\_B, fp32\_C$ )  $\triangleright$  Alg.12
7: FP32toU16 $\langle N/2, N/2 \rangle (c, fp32\_C)$   $\triangleright$  Alg.11

```

Algorithm 10 ParU16toFP16: Pre-processing elements for the matrix-vector product converting from U16 to FP16

Input: $N \times N$ -length polynomial in_1 and N -length polynomial in_2

Output: Matrix a_1, a_2, a_3 and vector b_1, b_2, b_3 in FP16 format

```

1:  $tid = thread\ ID$ 
2:  $bid = block\ ID$ 
// Launch  $N/2$  blocks and  $N/2$  threads in
// parallel
3: Algorithm 6 steps.

```

5 pre-processes polynomials A and B for the three TMVP multiplications in equation 2.

The arrangement of matrices and vectors is described in Algorithm 10. Similarly, line 6 executes Algorithm 12, which computes the matrix-vector product using Tensor-cores. Finally, line 7 post-processes the products and calculates the final result. Note that although the pre-processing steps for CUDA-cores (algorithms 5 and 6) and Tensor-cores (algorithms 5 and 10) are similar, the format of the output from Algorithms 6 and 10 differs. In CUDA-cores, the output is the same U16 format, whereas in Tensor-cores, the format changes to FP16.

Algorithm 12 shows the Tensor-cores polynomial convolution that computes all three multiplications in TMVP form. The matrix multiplication in Tensor-cores is performed as 16×16 having 32 threads in a warp. For larger matrices, multiple warps can be used to compute separate portions of the matrix. The results are then aggregated repeatedly to produce the final results. For example, to multiply a 32×32 matrix, four warps are launched in parallel to perform 16×16 matrix multiplication as shown in Figure 4. The other four warps compute the other half of the matrix in parallel. This means the process requires two iterations to perform 32×32 matrix multiplication. The final results are stored in parallel in Matrix C. However, for TMVP polynomial convolution in both Saber and Sable, $(128/16)^2$ warps, and $128/16$ iterations are required to perform the operation.

In Algorithm 12, matrix $a_1, a_2,$ and a_3 are comprised

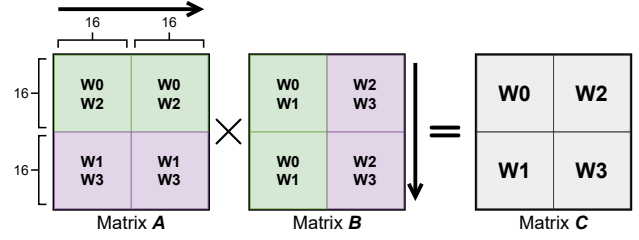


Fig. 4. Matrix multiplication in Tensor-cores of 32×32 having warps (w) running in parallel

Algorithm 11 FP32toU16: process polynomial coefficients from FP32 to U16 via three $n/2$ -dimensional TMVPs and modulo p

Input: $N/2 \times N/2$ matrix p_1, p_2 and p_3 with elements in FP32 format

Output: Matrix out of N -length degree, with elements in U16 format and modulo p

```

1:  $tid = thread\ ID$ 
2:  $bid = block\ ID$ 
// Launching  $N/2$  threads at maximum
3: if ( $tid < N/2$ ) then
4:    $out[bid + tid] += (int32\_t) (p_1[bid + tid] + p_2[bid +$ 
    $tid]) \% p$ 
5:    $out[bid + N/2 + tid] += (int32\_t) (p_1[bid + tid] -$ 
    $p_3[bid + tid]) \% p$ 
6: else
7:    $out[bid + tid] = 0$ 
8: end if

```

of public/private keys arranged and pre-processed in nega-cyclic form, and matrix b_1, b_2 and b_3 represent polynomial B . All matrices are stored in the global memory. Note that in this article, we use *fragment* to denote the temporary storage used to hold the matrices involved in Tensor-cores computations. First, Algorithm 12 initializes nine fragments: three for the 16×16 sub-matrices, three for sub-vectors, and three for collecting results of the multiplication of matrices and vector fragments (lines 1-9). The first multiplication iterates through matrix a_1 (row-major) and matrix b_1 (column-major) to multiply in parallel (lines 18-22). In each iteration, 16×16 sub-matrices are loaded from matrix a_1 and matrix b_1 (in global memory) for concurrent matrix multiplication. $(N/32)$ Each warp operates on separate regions of matrix a_1 and matrix b_1 . The collected results are transferred to matrix p_1 in global memory (line 24) in column-major form to ensure correctness. This is repeated for the other two multiplications (lines 26–40) and the outputs are stored in p_2 and p_3 .

Finally, referring to Algorithm 11, output matrix $p_1, p_2,$ and p_3 combine to get the final result (lines 4 and 5), and the format is converted from FP16 to U16. This whole process is done by using $N/2$ threads.

Algorithm 12 Tensor-cores: TMVP based parallel polynomial convolutions.

Input: $N/2 \times N/2$ -length matrices a_1, a_2, a_3 and $N/2$ -length vector b_1, b_2, b_3 , where N is a multiple of 16.

Output: $N/2 \times N/2$ -length matrix, p_1, p_2 , and p_3 holds the nega-cyclic convolution of distinct polynomials (a_1, b_1) , (a_2, b_2) , and (a_3, b_3) .

```

// 16 × 16 with precision FP16 initialization
of fragment (a1, b1), (a2, b2), & (a3, b3)
1: fragment < a1, 16, 16, 16, half, row_major > a1_frag
2: fragment < b1, 16, 16, 16, half, col_major > b1_frag
3: fragment < a2, 16, 16, 16, half, row_major > a2_frag
4: fragment < b2, 16, 16, 16, half, col_major > b2_frag
5: fragment < a3, 16, 16, 16, half, row_major > a3_frag
6: fragment < b3, 16, 16, 16, half, col_major > b3_frag
// 16 × 16 with precision FP32 initialization
of fragment C
7: fragment < accumulator, 16, 16, 16, float > c1_frag
8: fragment < accumulator, 16, 16, 16, float > c2_frag
9: fragment < accumulator, 16, 16, 16, float > c3_frag
// Compute the warp ID and indices
10: tid = thread ID
11: bid = block ID
12: blockDim = block dimension
13: id_warp = (bid × blockDim + tid)/32
14: row_idx = (id_warp%(N/32)) × 16
15: col_idx = (id_warp/(N/32)) × 16
16: acc_idx = row_idx + col_idx × N/2
17: for i from 0 to (N/32) do
18:   a1_id = row_idx × N/2 + i × 16
19:   b1_id = col_idx × N/2 + i × 16
20:   load_matrix_sync(a1_frag, a1 + a1_id, N/2)
21:   load_matrix_sync(b1_frag, b1 + b1_id, N/2)
22:   mma_sync(c1_frag, a1_frag, b1_frag, c1_frag)
23: end for
// Store c1_frag output in p1
24: store_matrix_sync(p1+acc_idx, c1_frag, N/2, col_major)
25: for i from 0 to (N/32) do
26:   a2_id = row_idx × N/2 + i × 16
27:   b2_id = col_idx × N/2 + i × 16
28:   load_matrix_sync(a2_frag, a2 + a2_id, N/2)
29:   load_matrix_sync(b2_frag, b2 + b2_id, N/2)
30:   mma_sync(c2_frag, a2_frag, b2_frag, c2_frag)
31: end for
// Store c2_frag output in p3
32: store_matrix_sync(p2+acc_idx, c2_frag, N/2, col_major)
33: for i from 0 to (N/32) do
34:   a3_id = row_idx × N/2 + i × 16
35:   b3_id = col_idx × N/2 + i × 16
36:   load_matrix_sync(a3_frag, a3 + a3_id, N/2)
37:   load_matrix_sync(b3_frag, b3 + b3_id, N/2)
38:   mma_sync(c3_frag, a3_frag, b3_frag, c3_frag)
39: end for
// Store c3_frag output in p3
40: store_matrix_sync(p3+acc_idx, c3_frag, N/2, col_major)

```

TABLE III
PERFORMANCE COMPARISON OF TMVP-BASED POLYNOMIAL
CONVOLUTION USING TENSOR-CORES AND CUDA-CORES

Batch size (K)	CUDA-cores	Tensor-cores
	Throughput (1000 multiplications per second)	
1	9.91	13.8
8	91.326	110.558
32	310.89	443.89
64	461.56	893.348
128	577.2	1844.34
256	738.497	3654.79
512	811.230	6740.39
1024	851.13	10861.83

IV. EXPERIMENT RESULTS AND DISCUSSION

This section presents a series of experiments to assess the efficacy of our proposed methodology. These experiments were conducted on a workstation equipped with a 2.10GHz Intel Core i7-12700F CPU with 16GB of RAM and an NVIDIA RTX3060 Ti GPU with a 1410 MHz frequency and 8 GB GDDR6 memory.

A. Performance of TMVP-2 polynomial convolution

In this section, we compare the performance of TMVP polynomial convolution using both CUDA-cores and Tensor-cores. We launched $(N/32)^2$ warps and N threads per block for polynomial convolution using Tensor-cores and CUDA-cores, respectively. Based on the results presented in Figure 5 and Table III, it is evident that the proposed TMVP polynomial convolution with Tensor-cores outperformed CUDA-cores. Although the difference is small at the initial batch sizes, the throughput on CUDA-cores starts to saturate when the batch size exceeds 64. The difference between Tensor-cores and CUDA-cores versions becomes significant as the batch size increases beyond 64. For instance, at a batch size of 1, Tensor-cores is only $1.39\times$ faster than CUDA-cores. However, this difference increases to $8.31\times$ and $12.76\times$ at batch sizes of 512 and 1024, respectively.

The low performance from CUDA-cores is due to the limited shared memory in the GPU and the large number of read/write operations required for polynomial convolution. As mentioned in Section III-B, the limited shared memory is insufficient to hold multiple copies of vectors. Consequently, to locate the precise element in the vector for matrix multiplication, the modulo operation must be employed. Nonetheless, this operation acts as a conditional statement for each thread, resulting in a reduction in performance. Secondly, as seen in Algorithm 7, we are conducting three TMVP multiplications in a single kernel. Table IV depicts the number of read/write operations executed in one CUDA-cores kernel. Overall, 82,432 reads and 49,920 writes were performed in one kernel to accomplish three TMVP polynomial convolutions. According to Table II, when using Saber and Sable, the value of $\ell = 3$. This means that in order to complete one polynomial convolution, $82,432 \times 3$ read operations and $49,920 \times 3$ write operations are required. However, the TMVP in Tensor-cores does not have the same memory limitations and is capable of

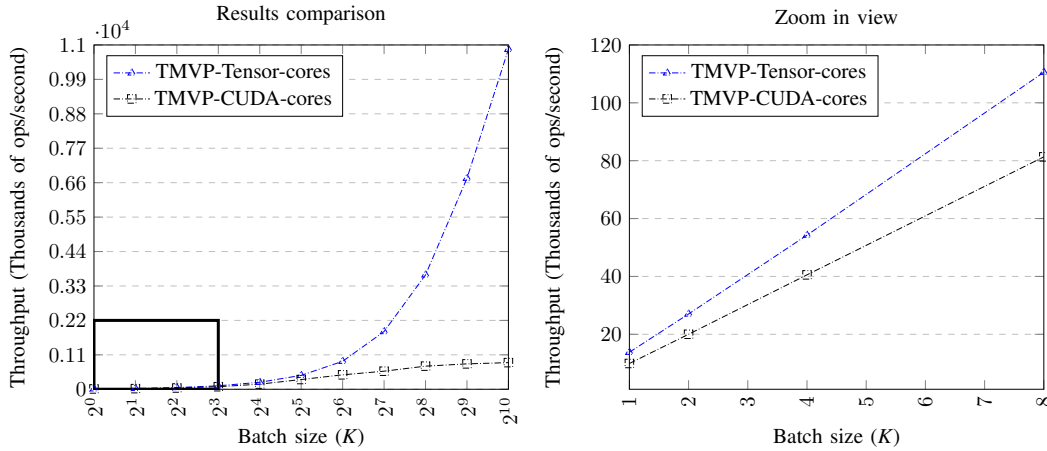


Fig. 5. Performance Comparison of TMVP-based polynomial convolution using Tensor-cores and CUDA-cores

TABLE IV
READ/WRITE OPERATIONS ON SHARED MEMORY FOR THE TMVP IN
CUDA-CORES

TMVP Multiplications	Tot. Read Elements	Tot. Write Elements	Tot. Read/Write Operations
P1	16640	16640	33280
P2	32896	16640	49536
P3	32896	16640	49536
Total Operations	82432	49920	132352

processing multiple copies of vectors simultaneously, resulting in high throughput compared to the TMVP in CUDA-cores.

B. Performance breakdown for TMVP Tensor-core and CUDA-core based implementations

Table V provides the performance breakdown of polynomial convolution in Saber and Sable by utilizing the proposed techniques on Tensor-cores and CUDA-cores. The analysis of execution times was conducted using a batch size of $K = 128$, with both CUDA-cores and Tensor-cores given a sufficient workload. In the tensor-cores version, organizing poly a into a nega-cyclic matrix takes up about 35% of the overall time. In contrast, pre-arrangement of poly A and poly B takes up 15% of the total time. Matrix-vector multiplication in Tensor-cores is the most time-consuming operation (about $\approx 37\%$), which is close to the nega-cyclic arrangement of poly a . Converting the format from FP32 to U16 and simultaneously performing reduction requires the least amount of time (about $\approx 12\%$). The performance breakdown in CUDA-cores shows that nega-cyclic rearrangement of poly a consumes only $\approx 11\%$ of the total time, whereas pre-arrangement of both polynomials only takes 5% of the time. Matrix-vector multiplication in CUDA-cores consumes the most time (about $\approx 79\%$). Post-arrangements of elements and performing modulo consume the least amount of time (nearly 4%).

Based on the above discussion, it becomes apparent that polynomial convolution using Tensor-cores requires less shorter multiplication time compared to CUDA-cores. This can be attributed to the superior capabilities provided by Tensor-

cores in executing matrix-vector multiplications more efficiently than CUDA-cores. Moreover, as elucidated in Section IV-A, the large number of read/write operations on shared memory required by CUDA-cores also needs more time for multiplication. This duration increases in proportion to larger batch sizes. In contrast, Tensor-cores do not use shared memory because most of the computations are performed directly in the registers.

C. Comparing KEX and KEM performance on a GPU

This section presents the KEX and KEM experiment results from Saber and Sable after implementing the TMVP techniques as proposed. The experiments take into account different batch sizes (K) and utilize two types of GPU: CUDA-cores and Tensor-cores. The KEX performance for both schemes is given in Table VI. Implementation of Saber encryption using TMVP on CUDA-cores and Tensor-cores yielded impressive results.

At a batch size of 16, Tensor-cores achieved 45,625 and 237,869 encryption and decryption operations per second, respectively, whereas CUDA-cores achieved only 35,358 and 179,921 operations per second. Notably, Tensor-cores was $1.2\times$ faster for encryption and $1.3\times$ faster for decryption than the CUDA-cores. We determined that increasing the batch size led to higher throughput for both schemes due to the increased workload in fully occupying a GPU. However, the difference in throughput between Tensor-cores and CUDA-cores also increased, with the highest throughput occurring at $K = 512$. In fact, Tensor-cores achieved 424,437 and 6,259,781 encryption and decryption operations per second, respectively, which were $2.6\times$ and $6.8\times$ faster than CUDA-cores. Similar results were observed with our implementation of Sable encryption and decryption. Initially, the difference between CUDA-cores and Tensor-cores encryption and decryption was small. At $K = 512$, Tensor-cores achieved 457,155 and 5,621,925 encryption and decryption operations per second, respectively, which were $2.7\times$ and $6.2\times$ faster than the CUDA-cores.

Table VII shows the throughput from Saber and Sable KEMs on a GPU using the TMVP on Tensor-cores and CUDA-cores. It is important to note that KEM is an extension

TABLE V
PERFORMANCE BREAKDOWN OF THE TMVP POLYNOMIAL CONVOLUTION USING TENSOR-CORES AND CUDA-CORES AT $K=128$

Operation	Tensor-cores		CUDA-cores	
	Time (μs)	%	Time (μs)	%
ParNegCyc (Poly $A \rightarrow$ Algorithm 5)	25.74	35.08	25.83	11.64
ParU16toFP16 (Pre-arrangement of Poly $A \& B \rightarrow$ Algorithm 10)	11.52	15.6	-	-
PreArr (Pre-arrangement of Poly $A \& B \rightarrow$ Algorithm 6)	-	-	11.82	5.33
Tensor-cores (Matrix-vector multiplication \rightarrow Algorithm 12)	27.18	37.04	-	-
CUDA-cores (Matrix-vector multiplication \rightarrow Algorithm 7)	-	-	175.29	79.04
FP32toU16 (Post arrangement \rightarrow Algorithm 11)	8.94	12.18	-	-
PostProcess (Post arrangement \rightarrow Algorithm 8)	-	-	8.82	4.00
Total	73.38	100	221.76	100

TABLE VI
COMPARING THE THROUGHPUT OF SABER AND SABLE KEX WITH THE TMVP CUDA-CORES AND TENSOR-CORES IMPLEMENTATIONS AT DIFFERENT BATCH SIZES

Batch size (K)	Saber				Sable			
	Throughput (encryptions/decryptions per second)							
	CUDA-cores		Tensor-cores		CUDA-cores		Tensor-cores	
	Encrypt	Decrypt	Encrypt	Decrypt	Encrypt	Decrypt	Encrypt	Decrypt
16	35358	179921	45625	237869	37838	170706	45183	235404
32	66401	338410	86821	498256	70546	316055	87412	457456
64	98068	505945	155436	957854	102838	483675	156678	896861
128	127218	647564	257583	1912960	129941	629029	263695	1757469
256	153852	824997	359680	3546473	161075	811112	374619	3218020
512	164670	916223	424437	6259781	170956	902628	457155	5621925

of KEX and involves additional hashing operations, which results in lower throughput compared to KEX. At batch size $K = 512$, the throughput of Saber on Tensor-cores technique was $2.0\times$ faster (encapsulation) and $2.37\times$ faster (decapsulation) than on CUDA-cores implementation. For Sable, the throughput for encapsulation and decapsulation was $1.9\times$ and $2.35\times$ higher than on CUDA-cores implementation, respectively.

A thorough analysis found that the matrix-vector multiplication on CUDA-cores is the most time-consuming, particularly when handling large batch sizes. This is primarily attributed to the fact that when K exceeds 64, CUDA-cores are fully loaded, where adding more workload (i.e., increasing the batch size) does not increase the throughput. To achieve optimal performance from GPU implementation, it is critical to utilize fast shared memory, but transferring between global and shared memory also have significant overhead. In contrast, Tensor-cores offer faster processing times owing to their accelerated matrix operations and smaller matrices available for multiplication in the TMVP. Although we require three TMVPs with 128×128 matrices (see Algorithm 12) instead of one 256×256 multiplication, the total number of operations executed for three TMVPs is still less than that in one 256×256 multiplication, resulting a faster polynomial convolution. Detailed explanation of these points can be found in sections III-A and III-C.

D. Comparison with State-of-the-Art implementations

1) Performance comparison with GPU implementations:

The graphs in Figure 6 show the performance comparison of proposed TMVP CUDA-cores and Tensor-cores polynomial

convolution with schoolbook polynomial convolution proposed by Hafeez et al. [26]. The results indicate that the schoolbook technique implemented on CUDA-cores (SB-CUDA) initially demonstrated impressive performance. However, as the batch size increased to exceed 64, the CUDA cores were fully utilized, and its performance began to saturate, leading to the saturation of the SB-CUDA throughput. As a result, the Tensor-cores implementation surpassed the CUDA-cores. It is worth noting that the throughput of the Tensor-cores continually increased as the batch size increased, even at a batch size greater than 64. This is the reason both Tensor-core approaches provide significantly higher throughput when the batch size is large, and both Tensor-core techniques exhibit similar performance until $K=256$. However, at $K \geq 512$, the proposed TMVP approach outperformed the schoolbook Tensor-cores (SB-TC) approach. It is important to note that the performance of the TMVP on CUDA-cores was even slower than SB-CUDA [26]. This shows that the TMVP may not always provide performance superior to the schoolbook approach because memory movement is critical in achieving good performance in GPU implementation.

Table VIII presents a throughput comparison of our proposed technique, with Schoolbook Tensor-cores implementation (SB-TC) [26] and dot-product instructions (DPSaber) [27]. Hafeez et al. [26] proposed SB-TC for Sable, and Lee et al. [27] proposed DPSaber for Saber. Note that Sable KEM is an improvement over Saber because it employs polynomial convolution for efficient inner product and matrix-vector multiplication calculations. DPSaber [27] incorporates dot-product instructions found in GPUs to implement Saber,

TABLE VII
COMPARING THE THROUGHPUT AT DIFFERENT BATCH SIZES FOR THE SABER AND SABLE KEM TMVPS IN CUDA-CORES AND TENSOR-CORES

Batch size (K)	Saber				Sable			
	Throughput (encaps/decaps per second)							
	CUDA-cores		Tensor-cores		CUDA-cores		Tensor-cores	
	Encaps	Decaps	Encaps	Decaps	Encaps	Decaps	Encaps	Decaps
16	25264	24919	29860	30628	25604	25697	29472	30376
32	46803	46572	55878	58268	46790	47221	53395	57336
64	74206	71411	101569	105983	73118	71648	96950	102769
128	102249	94271	171248	179163	97370	93989	156109	173205
256	123648	115786	229200	245263	121304	115829	212816	241560
512	133919	124261	267720	294020	130675	125340	250062	295061

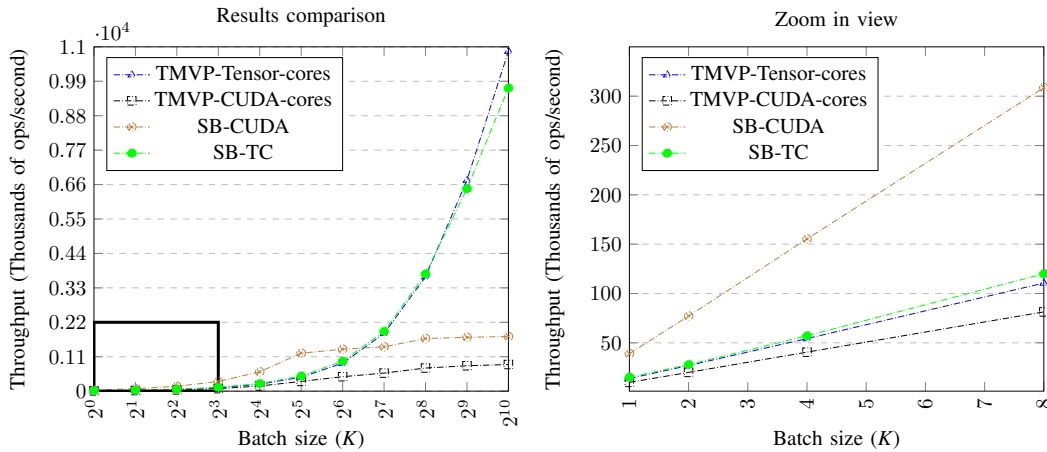


Fig. 6. Performance comparison of TMVP and Schoolbook based polynomial convolution using Tensor-cores and CUDA-cores

and SB-TC uses Tensor-cores for the schoolbook method for polynomial convolution in Sable. We conducted experiments on the same GPU used in the SB-TC [26] and directly adopted source code available in the public domain. Similarly, since the source code for DPSaber is open, we utilized that code and experimented on the same GPU for a fair comparison.

Table VIII provides insight into the performance of our proposed TMVP-TC version in comparison to SB-TC and DPSaber. Sp-up 1 denotes the ratio of TMVP-TC to SB-TC, while Sp-up 2 denotes the ratio of TMVP-TC to DPSaber. Referring to matrix-vector multiplication, our findings show that DPSaber performed better when $K \leq 64$. However, SB-TC achieved almost the same throughput as TMVP-TC. At $K \geq 128$, TMVP-TC outperformed DPSaber and achieved $4.24\times$ higher throughput at $K = 1024$. Similarly, in comparison to S-TC, TMVP-TC achieved $1.12\times$ higher throughput. Furthermore, our TMVP-TC achieved at least $3.63\times$ higher throughput than DPSaber for the inner product at $K = 1024$, but with SB-TC, we achieved $1.21\times$ higher throughput. These results demonstrate that our approach is more advantageous than SB-TC and DPSaber when the batch size is sufficiently large. This is due to the small matrices and fewer multiplication operations in the TMVP approach, as well as the higher instruction throughput on Tensor-cores in comparison to the dot-product instructions and the SB-TC approach.

2) Performance comparison with AVX2 implementations:

Table IX presents a comparative throughput analysis of inner

product and matrix-vector multiplication in Saber KEM using TMVP-TC against the AVX2 implementation proposed in Chung et al. [47].

In reference to inner product operation, AVX2 implementation shows better performance when $K \leq 64$ compared to the proposed TMVP-TC approach. However, TMVP-TC outperformed the AVX2 implementation at larger batch sizes and achieved $11.30\times$ higher throughput at batch size 1024 . Similarly, TMVP-TC outperforms AVX2 for matrix-vector multiplication even at $K = 64$ and shows $1.351\times$ higher throughput. This improvement increases as batch size increases, recording the highest throughput of $17.13\times$ faster at $K = 1024$. This highlights the efficiency of the proposed TMVP-TC method in handling large-scale operations, demonstrating its potential applicability in optimizing cryptographic computations in Saber KEM.

Table X shows the performance of TMVP-TC and AVX2 implementation for the Sable KEM. TMVP-TC outperforms AVX2 across various batch sizes in both encapsulation and decapsulation processes. For instance, at a batch size of 16, TMVP shows a throughput of 29,472 and 30,376 for encapsulation and decapsulation, respectively, compared to AVX2's performance of 27,368 and 29,213. The performance advantage of TMVP-TC becomes more pronounced as batch size increases. At a batch size of 512, TMVP achieves a throughput of 250,062 ($9.13\times$ faster) for encapsulation and 295,061 ($10.10\times$ faster) for decapsulation than AVX2.

TABLE VIII
PERFORMANCE COMPARISON OF TMVP ON TENSOR-CORES FOR INNER-PRODUCT AND MATRIX-VECTOR MULTIPLICATION IN THE SABER AND SABLE KEM VERSUS SCHOOLBOOK TENSOR-CORES [26] AND DPSABER [27] APPROACHES

Batch size (K)	Inner Product (thousand of operations per second)					Matrix-vector (thousand of operations per second)				
	TMVP-TC	SB-TC [26]	DPSaber [27]	Sp-up ¹	Sp-up ²	TMVP-TC	SB-TC [26]	DPSaber [27]	Sp-up ¹	Sp-up ²
64	893	957	1161	0.93	0.77	366	353	445	1.03	0.82
128	1844	1910	1926	0.96	0.96	762	711	734	1.07	1.07
256	3655	3746	2598	0.97	1.41	1495	1362	1001	1.09	1.49
512	6740	6465	2832	1.04	2.38	2795	2553	1034	1.09	2.70
1024	10861	9681	2991	1.21	3.63	4643	4144	1096	1.12	4.24

¹ TMVP-TC / SB-TC; ² TMVP-TC / DPSaber

TABLE IX
PERFORMANCE COMPARISON OF TMVP ON TENSOR-CORES FOR INNER-PRODUCT AND MATRIX-VECTOR MULTIPLICATION IN THE SABER KEM VERSUS [47] APPROACH

Batch size (K)	Inner Product (thousand of operations per second)			Matrix-vector (thousand of operations per second)		
	TMVP-TC	AVX2 [47]	Sp-up	TMVP-TC	AVX2 [47]	Sp-up
64	893		0.93	366		1.35
128	1844		1.92	762		2.81
256	3655	961	3.80	1495	271	5.52
512	6740		7.01	2795		10.31
1024	10861		11.30	4643		17.13

TABLE X
PERFORMANCE COMPARISON OF TMVP ON TENSOR-CORES AND AVX2 [7] IMPLEMENTATION FOR THE SABLE KEM

Batch size (K)	TMVP-TC		AVX2 [7]		Speed-up	
	Throughput (encapsulation/decapsulation per second)					
	Encap	Decap	Encap	Decap	Encap	Decap
16	29472	30376			1.08	1.04
32	53395	57336			1.95	1.96
64	96950	102769			3.54	3.51
128	156109	173205	27368	29213	5.70	5.92
256	212816	241560			7.78	8.27
512	250062	295061			9.13	10.10

3) Performance comparison with Kyber implementations:

To benchmark our proposed methods against Kyber (a finalist in the NIST PQC standardization process) provides a relevant and stringent test of our proposed method's performance and robustness in the context of advanced cryptographic standards. Table XI provides the performance comparison of the KEX implementation of Saber and Sable with the NIST standardized implementation of Kyber on GPU by Lee et al. [20].

As the batch size increases for Saber and Sable, so does the throughput. At a batch size of $K = 16$, Saber has an encryption throughput of 45,625 and a decryption throughput of 237,869, while Sable shows similar results with 45,183 and 235,404 for encryption and decryption respectively. Kyber, on the other hand, outperforms both with encryption and decryption throughputs of 98,502 and 361,271, respectively.

However, as the batch size increases, the throughput gap between Saber, Sable, and Kyber decreases. At a batch size 512, Saber and Sable reach their peak throughputs. Meanwhile, Kyber maintains its lead in encryption with a throughput of 1,002,645 compared to 424,437 and 457,155, respectively,

but Saber and Sable outperform Kyber in decryption with a throughput of 6,259,781 ($1.80\times$ faster) and 5,621,925 ($1.66\times$ faster) compared to the 3,393,425.

E. Application to the IoT Communication

The proposed TMVP technique introduces significant theoretical and engineering innovations for optimizing the polynomial convolution during the implementation of PQC algorithms. By leveraging TMVP, the computational complexity of polynomial convolution is reduced, leading to improved efficiency and performance. This approach not only streamlines the algorithmic process, but also minimizes storage requirements, making it well-suited for resource-constrained environments like IoT devices.

Moreover, gateway devices and servers typically handle the bulk of the data traffic and require high-throughput solutions. Adopting the proposed tensor-cores based TMVP technique for PQC algorithms offers significant enhancements in throughput, enabling high-speed encapsulation and decapsulation operations crucial for IoT applications. For instance, our Sable implementation achieves impressive rates of 250,062 encapsulations/s and 295,061 decapsulations/s on an RTX 3060Ti GPU. This enhanced throughput is well-suited for IoT scenarios, ensuring secure, rapid, and resource-efficient cryptographic operations across the IoT ecosystem.

V. CONCLUSION

Our research demonstrated the effectiveness of parallel TMVP computations utilizing Tensor-cores and CUDA-cores in accelerating the execution of KEX and KEM algorithms. By applying this technique to the post-quantum KEMs (Saber and Sable), we achieved significant improvements in system performance where high throughput is required, especially for

TABLE XI
PERFORMANCE COMPARISON OF KEX (BOTH SABER AND SABLE) WITH NIST STANDARDIZATION KYBER IMPLEMENTATION [20]

Batch size (K)	Saber (TMVP-TC)		Sable (TMVP-TC)		Kyber [20]		Sp-up ¹		Sp-up ²	
	Throughput (encryption/decryption per second)									
	Encrypt	Decrypt	Encrypt	Decrypt	Encrypt	Decrypt	Encrypt	Decrypt	Encrypt	Decrypt
16	45625	237869	45183	235404	98502	361271	0.46	0.66	0.46	0.65
32	86821	498256	87412	457456	161838	677966	0.54	0.73	0.54	0.67
64	155436	957854	156678	896861	297044	1212856	0.52	0.78	0.53	0.74
128	257583	1912960	263695	1757469	596836	1806684	0.43	1.06	0.44	0.97
256	359680	3546473	374619	3218020	838222	2633311	0.43	1.34	0.45	1.22
512	424437	6259781	457155	5621925	1002645	3393425	0.42	1.80	0.45	1.66

¹ Saber (TMVP-TC) / Kyber [20]; ² Sable (TMVP-TC) / Kyber [20]

IoT applications. In the case of Sable, our proposed Tensor-cores implementation outperformed traditional CUDA-cores implementations in terms of encryption and decryption speeds. Specifically, we achieved a minimum of $1.1\times$ faster encryption and $1.07\times$ faster decryption. Moreover, our approach demonstrated $1.7\times$ higher throughput for encryption and an impressive $3.1\times$ higher throughput for decryption in KEX operations.

In the future, we plan to work on overcoming limitations related to GPU shared memory and exploring the potential for parallel multiplications to improve throughput further. This could be achieved by caching some content onto the registers to lower the shared memory usage. Furthermore, we aim to optimize the arrangement and processing of polynomial coefficients to minimize computation time and boost the efficiency of other lattice-based cryptographic schemes that do not directly support NTT operations on GPUs. This may involve exploring advanced GPU architectures, refining algorithms to minimize read/write operations on shared memory, and experimenting with larger batch sizes to capitalize on GPU capacities for higher throughput in cryptographic applications.

ACKNOWLEDGMENT

This work was supported by the Gachon University research fund under Grant GCU-202304050001 and was conducted by funding from The Circle Foundation (Republic of Korea) for 1 year since December 2023 as Quantum Security Research Center selected as the 2023 The Circle Foundation Innovative Science Technology Center under Grant 2023 TCF Innovative Science Project-05.

REFERENCES

- [1] I. T. L. Computer Security Division, "Post-quantum cryptography standardization - post-quantum cryptography: CSRC." [Online]. Available: <https://csrc.nist.gov/Projects/Post-Quantum-Cryptography/Post-Quantum-Cryptography-Standardization>
- [2] J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, P. Schwabe, G. Seiler, and D. Stehlé, "CRYSTALS-Kyber: a CCA-secure module-lattice-based KEM," in *2018 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2018, pp. 353–367.
- [3] V. Lyubashevsky, L. Ducas, E. Kiltz, T. Lepoint, P. Schwabe, G. Seiler, D. Stehlé, and S. Bai, "Crystals-Dilithium," *Algorithm Specifications and Supporting Documentation*, 2020.
- [4] P.-A. Fouque, J. Hoffstein, P. Kirchner, V. Lyubashevsky, T. Pornin, T. Prest, T. Ricosset, G. Seiler, W. Whyte, Z. Zhang *et al.*, "Falcon: Fast-Fourier lattice-based compact signatures over NTRU," *Submission to the NIST's post-quantum cryptography standardization process*, vol. 36, no. 5, pp. 1–75, 2018.
- [5] J.-P. Aumasson, D. J. Bernstein, W. Beullens, C. Dobraunig, M. Eichlseder, S. Fluhrer, S.-L. Gazdag, A. Hülsing, P. Kampanakis, S. Kölbl *et al.*, "SPHINCS," 2019.
- [6] P.-Q. Cryptography, "Round 2 submissions," *Electronic resource. Access mode: https://csrc.nist.gov/Projects/post-quantum-cryptography/round-2-submissions*, 2021.
- [7] J. M. Bermudo Mera, A. Karmakar, S. Kundu, and I. Verbauwhede, "Scabbard: a suite of efficient learning with rounding key-encapsulation mechanisms," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, vol. 2021, no. 4, p. 474–509, Aug. 2021.
- [8] J.-P. D'Anvers, A. Karmakar, S. Sinha Roy, and F. Vercauteren, "Saber: Module-LWR based key exchange, CPA-secure encryption and CCA-secure KEM," in *Progress in Cryptology—AFRICACRYPT 2018: 10th International Conference on Cryptology in Africa, Marrakesh, Morocco, May 7–9, 2018, Proceedings 10*. Springer, 2018, pp. 282–305.
- [9] KpqC, "Korean pqc competition," 2022, [Online; accessed 30-June-2023]. [Online]. Available: <https://www.kpqc.or.kr/competition.html>
- [10] Z. Liang, B. Fang, J. Zheng, and Y. Zhao, "Compact and Efficient KEMs over NTRU Lattices," *Cryptology ePrint Archive*, 2022.
- [11] C. Chen, O. Danba, J. Hoffstein, A. Hülsing, J. Rijneveld, J. M. Schanck, P. Schwabe, W. Whyte, and Z. Zhang, "Algorithm specifications and supporting documentation," *Brown University and Onboard security company, Wilmington USA*, 2019.
- [12] J. Cho, Y. Lee, Z. Koo, J.-S. No, and Y.-S. Kim, "Improving Key Size and Bit-Security of Modified pqsigRM," in *2022 13th International Conference on Information and Communication Technology Convergence (ICTC)*. IEEE, 2022, pp. 1463–1467.
- [13] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Mathematics of computation*, vol. 19, no. 90, pp. 297–301, 1965.
- [14] C.-Y. Lee, P. K. Meher, and W.-Y. Lee, "Subquadratic space complexity digit-serial multiplier over binary extension fields using Toom-Cook algorithm," in *2014 International Symposium on Integrated Circuits (ISIC)*. IEEE, 2014, pp. 176–179.
- [15] Z.-Y. Wong, D. C.-K. Wong, W.-K. Lee, K.-M. Mok, W.-S. Yap, and A. Khalid, "KaratSaber: New Speed Records for Saber Polynomial Multiplication using Efficient Karatsuba FPGA Architecture," *IEEE Transactions on Computers*, 2023.
- [16] I. K. Paksoy and M. Cenk, "TMVP-based multiplication for polynomial quotient rings and application to saber on arm Cortex-M4," *Cryptology ePrint Archive*, 2020.
- [17] —, "Faster NTRU on ARM Cortex-M4 with TMVP-based multiplication," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 69, no. 10, pp. 4083–4092, 2022.
- [18] K. Townsend, "Solving the Quantum Decryption 'Harvest Now, Decrypt Later' Problem." [Online]. Available: <https://www.securityweek.com/solving-quantum-decryption->
- [19] N. Gupta, A. Jati, A. K. Chauhan, and A. Chattopadhyay, "PQC acceleration using GPUs: Frodokem, NewHope, and KYBER," *IEEE Transactions on Parallel and Distributed Systems*, vol. 32, no. 3, pp. 575–586, 2020.
- [20] W.-K. Lee and S. O. Hwang, "High throughput implementation of post-quantum key encapsulation and decapsulation on GPU for Internet of Things applications," *IEEE Transactions on Services Computing*, vol. 15, no. 6, pp. 3275–3288, 2021.
- [21] V. L. R. d. Costa, J. López, and M. V. Ribeiro, "A system-on-a-chip implementation of a post-quantum cryptography scheme for smart meter data communications," *Sensors*, vol. 22, no. 19, p. 7214, 2022.

- [22] M. Adeli, N. Bagheri, H. R. Maimani, S. Kumari, and J. J. Rodrigues, "A post-quantum compliant authentication scheme for iot healthcare systems," *IEEE Internet of Things Journal*, 2023.
- [23] S. Paul, P. Scheible, and F. Wiemer, "Towards post-quantum security for cyber-physical systems: Integrating pqc into industrial m2m communication 1," *Journal of Computer Security*, vol. 30, no. 4, pp. 623–653, 2022.
- [24] NVIDIA, "Jetson orin," 2023. [Online; accessed 31-October-2023]. [Online]. Available: <https://www.nvidia.com/en-us/autonomous-machines/embedded-systems/jetson-orin/>
- [25] Hardkernel, "Odroid-m1s," 2023. [Online; accessed 31-October-2023]. [Online]. Available: <https://www.hardkernel.com/shop/odroid-m1s-with-8gbyte-ram/>
- [26] M. A. Hafeez, W.-K. Lee, A. Karmakar, and S. O. Hwang, "High Throughput Acceleration of Scabbard Key Exchange and Key Encapsulation Mechanism Using Tensor Core on GPU for IoT Applications," *IEEE Internet of Things Journal*, 2023.
- [27] W.-K. Lee, H. Seo, S. O. Hwang, R. Achar, A. Karmakar, and J. M. B. Mera, "DPCrypto: Acceleration of Post-Quantum Cryptography Using Dot-Product Instructions on GPUs," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 69, no. 9, pp. 3591–3604, 2022.
- [28] H. Fan and M. A. Hasan, "A new approach to subquadratic space complexity parallel multipliers for extended binary fields," *IEEE Transactions on Computers*, vol. 56, no. 2, pp. 224–233, 2007.
- [29] M. A. Hasan, N. Meloni, A. H. Namin, and C. Negre, "Block recombination approach for subquadratic space complexity binary field multiplication based on Toeplitz matrix-vector product," *IEEE Transactions on Computers*, vol. 61, no. 2, pp. 151–163, 2010.
- [30] M. A. Hasan and C. Negre, "Multiway splitting method for Toeplitz matrix-vector product," *IEEE Transactions on Computers*, vol. 62, no. 7, pp. 1467–1471, 2012.
- [31] S. Ali and M. Cenk, "Faster residue multiplication modulo 521-bit Mersenne prime and an application to ECC," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 65, no. 8, pp. 2477–2490, 2018.
- [32] H. K. Taşkın and M. Cenk, "Speeding up curve25519 using Toeplitz matrix-vector multiplication," in *Proceedings of the Fifth Workshop on Cryptography and Security in Computing Systems*, 2018, pp. 1–6.
- [33] S. Winograd, *Arithmetic complexity of computations*. Siam, 1980, vol. 33.
- [34] J.-S. Pan, C.-Y. Lee, A. Sghaier, M. Zeghid, and J. Xie, "Novel systolization of subquadratic space complexity multipliers based on Toeplitz matrix-vector product approach," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 27, no. 7, pp. 1614–1622, 2019.
- [35] J. H. Cheon, H. Choe, D. Hong, and M. Yi, "SMAUG: Pushing Lattice-based Key Encapsulation Mechanisms to the Limits," *Cryptology ePrint Archive*, 2023.
- [36] A. Banerjee, C. Peikert, and A. Rosen, "Pseudorandom functions and lattices," in *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 2012, pp. 719–737.
- [37] W.-K. Lee, H. Seo, Z. Zhang, and S. O. Hwang, "Tensorcrypto: High throughput acceleration of lattice-based cryptography using tensor core on gpu," *IEEE Access*, vol. 10, pp. 20616–20632, 2022.
- [38] J.-C. See, H.-F. Ng, H.-K. Tan, J.-J. Chang, K.-M. Mok, W.-K. Lee, and C.-Y. Lin, "Cryptensor: A resource-shared co-processor to accelerate convolutional neural network and polynomial convolution," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, 2023.
- [39] Y. Gao, J. Xu, and H. Wang, "CuNH: Efficient GPU implementations of post-quantum KEM NewHope," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, no. 3, pp. 551–568, 2021.
- [40] S. Sun, R. Zhang, and H. Ma, "Efficient parallelism of post-quantum signature scheme SPHINCS," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 11, pp. 2542–2555, 2020.
- [41] W. Dai, B. Sunar, J. Schanck, W. Whyte, and Z. Zhang, "NTRU modular lattice signature scheme on CUDA GPUs," in *2016 International Conference on High Performance Computing & Simulation (HPCS)*. IEEE, 2016, pp. 501–508.
- [42] R. Avanzi, J. Bos, L. Ducas, E. Kiltz, T. Lepoint, V. Lyubashevsky, J. M. Schanck, P. Schwabe, G. Seiler, and D. Stehlé, "Crystals-kyber," *NIST, Tech. Rep.*, 2017.
- [43] V. B. Dang, K. Mohajerani, and K. Gaj, "High-speed hardware architectures and FPGA benchmarking of crystals-KYBER, NTRU, and Saber," *IEEE Transactions on Computers*, 2022.
- [44] Z. Chen, Y. Ma, T. Chen, J. Lin, and J. Jing, "Towards efficient Kyber on FPGAs: A processor for vector of polynomials," in *2020 25th Asia and South Pacific Design Automation Conference (ASP-DAC)*. IEEE, 2020, pp. 247–252.
- [45] G. Xin, J. Han, T. Yin, Y. Zhou, J. Yang, X. Cheng, and X. Zeng, "VPQC: A domain-specific vector processor for post-quantum cryptography based on RISC-V architecture," *IEEE transactions on circuits and systems I: regular papers*, vol. 67, no. 8, pp. 2672–2684, 2020.
- [46] F. Farahmand, V. B. Dang, M. Andrzejczak, and K. Gaj, "Implementing and benchmarking seven round 2 lattice-based key encapsulation mechanisms using a software/hardware codesign approach," in *Second PQC Standardization Conference*, 2019.
- [47] C.-M. M. Chung, V. Hwang, M. J. Kannwischer, G. Seiler, C.-J. Shih, and B.-Y. Yang, "Ntt multiplication for ntt-unfriendly rings: New speed records for saber and ntru on cortex-m4 and avx2," *IACR Transactions on Cryptographic Hardware and Embedded Systems*, pp. 159–188, 2021.



Muhammad Asfand Hafeez received a B.S. degree in electrical engineering from the University of Management and Technology in 2021. He is currently pursuing his master's degree in IT convergence engineering at Gachon University, South Korea. His research pursuits center on cryptography, GPU computing, deep learning, and hardware implementations.



Wai-Kong Lee received a B.Eng. in electronics and an M.Eng.Sc. from Multimedia University, Malaysia in 2006 and 2009, respectively. He received a Ph.D. in engineering from Universiti Tunku Abdul Rahman, Malaysia in 2018. Prior to joining academia, he worked in several multi-national companies including Agilent Technologies (Malaysia) as an R&D engineer. His research interests include cryptography, numerical algorithms, GPU computing, the Internet of Things, and energy harvesting. He is currently a post-doctoral researcher at Gachon University, South

Korea.



Angshuman Karmakar received the B.E. degree in computer science and engineering from Jadavpur University, Kolkata, India, in 2010 the M.Tech. degree in computer science and engineering from the Indian Institute of Technology, Kharagpur, India, in 2012 and the Ph.D. degree from Katholieke Universiteit Leuven (KU Leuven), Belgium, in 2020 for his dissertation titled "Design and Implementation Aspects of Post-Quantum Cryptography." He is one of the primary designers of the post-quantum Saber KEM scheme which is one of the finalists in the

NIST's post-quantum standardization procedure. He is currently working as an assistant professor at the Indian Institute of Technology, Kanpur, in India, and as a free researcher at COSIC, KU Leuven, Belgium. Earlier he was an FWO Post-Doctoral Fellow with the COSIC Research Group, KU Leuven. His research interest spans different aspects of lattice-based post-quantum cryptography and computation on encrypted data.



Seoung Oun Hwang received a B.S. degree in mathematics from Seoul National University, in 1993, the M.S. degree in information and communications engineering from the Pohang University of Science and Technology, in 1998, and a Ph.D. degree in computer science from the Korea Advanced Institute of Science and Technology, South Korea. He worked as a Software Engineer with LG-CNS Systems, Inc., from 1994 to 1996. He also worked as a Senior Researcher with the Electronics and Telecommunications Research Institute (ETRI), from 1998 to

2007. He worked as a Professor at the Department of Software and Communications Engineering, Hongik University, from 2008 to 2019. He is currently a Professor at the Department of Computer Engineering, at Gachon University. He is also an Editor of the ETRI Journal. His research interests include cryptography, cybersecurity, and artificial intelligence