

Low-Latency Linear Transformations with Small Key Transmission for Private Neural Network on Homomorphic Encryption*

Byeong-Seo Min^a, Joon-Woo Lee^{b,*}

^aDepartment of Electrical Engineering, Pohang University of Science and Technology (POSTECH), 77 Cheongam-ro, Nam-gu, Pohang, 37673, Gyeongsangbuk-do, South Korea

^bDepartment of Computer Science and Engineering, Chung-Ang University, 84 Heukseok-ro, Dongjak-gu, Seoul, 06974, South Korea

ARTICLE INFO

Keywords:

Privacy-Preserving Machine Learning
Fully Homomorphic Encryption
Cheon-Kim-Kim-Song (CKKS) schemes
Convolution

ABSTRACT

The widespread adoption of Machine Learning as a Service has made data privacy a critical challenge. Homomorphic Encryption (HE) offers a robust solution by enabling computations on encrypted data, yet it faces practical hurdles such as high computational latency and large key sizes. Consequently, designing efficient HE-based AI models is an active area of research. A primary focus is the HE implementation of convolution, a fundamental operation in CNNs and increasingly used in modern architectures such as transformers and state-space models. To realize HE convolution, various packing strategies—defining the data layout within a ciphertext—have been introduced, leading to state-of-the-art methods such as multiplexed parallel convolution (MPConv), which is implemented with multiplexed parallel packing. In this paper, we propose *Rotation-Optimized Multiplexed Parallel Convolution* (RO-MPConv), which enhances usability by reducing the number of rotation operations and rotation keys—major contributors to latency and key sizes in MPConv. Additionally, we introduce a *small-level key system* to further decrease key sizes and a novel *parallel baby-step giant-step matrix-vector multiplication*, which reduces rotations for packing strategies with multiple identical data entries, including multiplexed parallel packing. We conducted all experiments using the Lattigo library’s CKKS HE scheme. Experimental results demonstrate that RO-MPConv achieves up to an 81% reduction in latency and a 29× reduction in rotation key size compared to MPConv. Furthermore, integrating RO-MPConv into existing state-of-the-art models improved their total latency by up to 26%, and our proposed matrix-vector multiplication method reduced latency by 69%. Our code is fully available from <https://github.com/byeongseomin51/RO-MPConv.git>.

1. Introduction

As Machine Learning as a Service (MLaaS) becomes increasingly integral to various industries, the protection of client data has emerged as a critical challenge. This has catalyzed research into Privacy-Preserving Machine Learning (PPML), which aims to reconcile the utility of MLaaS with the fundamental right to data privacy.

A comprehensive survey by [1] identifies Homomorphic Encryption (HE), Differential Privacy (DP), and Secure Multi-Party Computation (MPC) as primary technologies for achieving PPML, each with distinct trade-offs. For instance, DP ensures privacy by adding statistical noise, which may come at the cost of model accuracy. MPC offers robust security for collaborative computation but often incurs high communication costs, posing a significant drawback in client-server scenarios. In contrast, Homomorphic Encryption presents a compelling alternative by enabling direct computation on encrypted data. This approach preserves

data utility, thereby maintaining high model performance, and provides strong cryptographic privacy guarantees, establishing it as a robust foundation for secure MLaaS.

In HE-based PPML, deep neural networks are typically implemented using Fully Homomorphic Encryption (FHE), as its bootstrapping functionality is essential for handling the substantial number of sequential operations required by deep architectures. Among the various FHE schemes, the Residue Number System variant of the Cheon-Kim-Kim-Song (RNS-CKKS) scheme [2], [3] is widely adopted for its native support for real-number arithmetic, a critical requirement for neural network computations. In standard FHE-based PPML, the inference process begins with a client encrypting their data and sending it with evaluation keys to a server. The server performs FHE model inference on the encrypted data and returns an encrypted result, which only the client can decrypt. Despite its security strengths, FHE-based PPML faces practical challenges, namely the significant computational latency of FHE inference and the large size of the required keys. Consequently, enhancing the efficiency of FHE-based PPML remains a central focus of current research.

Research on improving HE-based models can be broadly classified into two main avenues: the optimization of linear and non-linear operations. Homomorphic encryption natively supports linear operations such as addition, multiplication, and rotation, allowing the results of these computations to be nearly identical to those performed on plaintext. However, a key challenge arises from the data layout within

*This research was supported by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No.RS-2024-00398360, Development of a user-friendly and efficiency-optimized real-time homomorphic statistical analysis processing platform & No.RS-2024-00399401, Development of post-quantum security infrastructure transition and comprehensive quantum security verification technology).

*Corresponding author.

 minbyeongseo@postech.ac.kr (B. Min); jwlee2815@cau.ac.kr (J. Lee)

ORCID(s): 0009-0007-3944-8653 (B. Min); 0000-0002-4125-6331 (J. Lee)

a ciphertext, which stores data in a one-dimensional array of slots. The performance of multi-dimensional operations, therefore, critically depends on the packing strategy—the method used to map high-dimensional data onto this one-dimensional structure. The choice of this strategy directly dictates the efficiency of homomorphic linear operations, which is measured by the number of computationally expensive operations required, such as rotations and ciphertext-ciphertext multiplications, and the consumption of the multiplicative level.

Conversely, non-linear operations, such as activation functions, are generally handled by first approximating them with polynomials, which can then be evaluated using HE operations. For polynomial approximations, the packing strategy has a less direct impact on performance. Instead, a fundamental trade-off exists: higher-degree polynomials yield better accuracy but consume more multiplicative levels. The rapid consumption of the multiplicative level, in turn, increases the frequency of bootstrapping—the most computationally intensive operation in FHE—which is required to refresh the ciphertext’s capacity for further multiplications.

In this paper, we focus on the optimization of linear operations, with a specific emphasis on the convolution operation. Convolution is a cornerstone of Convolutional Neural Networks (CNNs) and is increasingly being integrated into modern Artificial Intelligence (AI) architectures such as Transformers [4, 5, 6] and State-Space Models (SSMs) [7]. To implement HE-based convolution, various packing strategies have been proposed. Among these, multiplexed parallel packing [8] is recognized as a state-of-the-art approach that enhances ciphertext efficiency by compactly filling the empty slots that arise from strided convolutions and by arranging multiple small input data blocks contiguously. This packing strategy is the foundation for multiplexed parallel convolution (MPConv), a state-of-the-art FHE convolution algorithm introduced in [8] to implement ResNet inference on multiplexed parallel packed data (hereafter MP-CNN) for the CIFAR-10 dataset. The significance of MPConv is underscored by its adoption in subsequent research, such as in AutoFHE [9] or CryptoFace [10].

However, MPConv suffers from two significant drawbacks: a high number of required rotation operations and a correspondingly large set of rotation keys. [8] mitigated this issue by strategically designing the network architecture to execute MPConv operations at the lowest possible multiplicative levels, where rotations are fastest. While this design is effective for MP-CNN, it is not a general solution. As research into non-linear function approximation advances, the multiplicative level consumption can change, and modifications to the CNN architecture may make it difficult or impossible to constrain MPConv operations to these low levels.

To overcome these limitations, this paper makes the following key contributions:

1. We propose *Rotation-Optimized Multiplexed Parallel Convolution* (RO-MPConv), an algorithm that reconstructs MPConv to reduce the number of rotations

and the volume of rotation keys without other performance degradation. This approach directly addresses the fundamental challenges of latency and key volume in FHE-based PPML. Our experiments show that RO-MPConv reduces latency by up to **56%** and rotation key size by up to **95%** compared to the original MPConv.

2. To provide greater flexibility for designing privacy-preserving convolutional layers, we introduce a method that enables a trade-off between consuming additional multiplicative levels and reducing the number of rotation operations. By strategically consuming one or two additional levels, our RO-MPConv variants can achieve further latency reductions of up to **79%** and **81%**, respectively. Crucially, we demonstrate that simply replacing MPConv with our RO-MPConv in existing state-of-the-art models improves their **total inference latency by up to 26%**.
3. We further propose a *small-level key system* that leverages the fact that convolution operations are typically confined to a specific range of multiplicative levels. When combined with the existing hierarchical rotation key system [11], this method reduced the rotation key volume by approximately 29-fold compared to MP-CNN.

As an additional contribution, we introduce parallel Baby-Step Giant-Step (BSGS) matrix-vector multiplication to accelerate operations like fully connected layers. By requiring fewer rotations than the standard BSGS algorithm with the diagonal method in [12], our algorithm achieves a latency reduction of up to 69%. It is especially efficient in settings where a ciphertext holds multiple identical data entries sequentially, a common feature in various packing strategies, including multiplexed parallel packing.

2. Related Works

The efficiency of deep neural networks over homomorphic encryption is fundamentally tied to the chosen packing strategy, which dictates how multi-dimensional data is mapped onto one-dimensional ciphertexts. The optimal packing strategy is often problem-dependent, varying with the model architecture and its specific requirements. Among the existing approaches, multiplexed parallel packing [17] has emerged as a state-of-the-art strategy for implementing HE-based CNNs and their variants. Its ability to support generalized strided convolutions and its modular, block-like applicability make it a versatile solution. Consequently, it is widely adopted in research aiming to enhance the practicality of HE-based models.

The broad applicability of MPConv is evident in its use across various advanced research domains. For instance, AutoFHE [9] and LPFHE [13] employed MPConv while exploring mixed-degree CNNs, which utilize polynomials of varying degrees for non-linear approximations across the network instead of a single, fixed degree. CryptoFace [10] implemented a private face recognition system using

a Patch CNN architecture, employing MPConv to perform the convolution operations. Research into HE-specific hardware accelerators or compilers [14, 15, 16, 17] frequently uses MP-CNN as a primary ResNet baseline for performance evaluation, underscoring its status as a benchmark in the field. The method's versatility has also been demonstrated in hybrid-packed CNNs [18] and Graph Convolutional Networks (GCNs) [19]. Furthermore, it is orthogonal to techniques like group convolution [20], summed area table (SAT) in [21], suggesting that these methods can be combined for even greater performance gains.

While multiplexed parallel packing is highly effective, it is important to acknowledge the diverse landscape of alternative packing and encoding schemes developed to address specific challenges in HE-based CNNs. In this context, packing is the strategy for arranging high-dimensional tensor data into a one-dimensional vector, whereas encoding is the process of transforming this vector into a plaintext polynomial that the encryption system can process. The choice of encoding directly influences the packing strategy. The most conventional approach is slot encoding, which is the method utilized by the MP-CNN baseline.

Other notable approaches are as follows: Falcon [22] employs frequency-domain packing with spectral encoding, while Gazelle [23] uses a packed Single-Input Single-Output (SISO) packing method. Some works focus primarily on novel encoding schemes, such as the optimal group encoding in CrossNet [24], coefficient encoding in [25], and the Coefficients-in-Slot (CinS) encoding in Neujeans [26]. Others, like Orion [27] and Chet [28], introduce compiler-level optimizations, such as Toeplitz matrix-based packing and automatic data layout selection, respectively.

However, each of these specialized methods entails trade-offs that limit their general applicability. For instance, Falcon's frequency-domain approach requires computationally expensive Fourier transforms for non-convolutional layers. The methods in CrossNet and Gazelle are designed for interactive settings that require frequent client communication, unlike the non-interactive scenarios we consider. The techniques in Chet and Orion are often tightly coupled with their specific compilers, reducing their flexibility for custom modifications. Furthermore, coefficient-based encodings like that of [25] sacrifice SIMD efficiency, and joint optimization strategies like the convolution-bootstrapping fusion in Neujeans [26] are ineffective when convolution and bootstrapping operations are not performed sequentially, such as in mixed-degree models.

Ultimately, no single packing or encoding strategy has proven universally dominant for all HE applications. Each approach presents a unique set of advantages and trade-offs. However, because multiplexed parallel packing is model-agnostic and highly applicable to a diverse range of architectures and practical use cases, it remains one of the most robust and widely used methods. Therefore, optimizing its core algorithm, MPConv, is a critical step toward bridging the gap between theoretical HE-based models and their real-world deployment. By enhancing the efficiency of

this foundational algorithm, our work contributes to making complex, privacy-preserving AI systems more practical and accessible.

3. Preliminaries

3.1. RNS-CKKS Fully Homomorphic Encryption

The residue number system variant Cheon-Kim-Kim-Song (RNS-CKKS) is a Fully Homomorphic Encryption (FHE) scheme that enables real number operations on encrypted data. In RNS-CKKS, all ciphertexts consist of one-dimensional complex number data. Specifically, they are in the form of $(b, a) \in R_Q^2$, where Q is a product of some prime numbers and $R_Q = \mathbb{Z}_Q[X]/\langle X^N + 1 \rangle$. A single ciphertext contains $N/2$ slots, each slot consisting of a single complex number. In this paper, $N/2$ is represented as n_t . All plaintext also has a one-dimensional shape, and for all one-dimensional data including plaintext in this paper, assuming there is a one-dimensional data A , $A[i]$ denotes the i -th data of A (where i starts from 0). Similarly, this applies equally to higher-dimensional data. For instance, in two-dimensional data \bar{A} , the value at the i -th row and j -th column is denoted as $\bar{A}[i][j]$.

In RNS-CKKS, there are three main homomorphic operations: addition, multiplication, and rotation. The addition and multiplication operations correspond to addition and multiplication between slots at the same positions within the ciphertexts, respectively. Rotation involves cyclic shifting of the values stored in each slot. In our paper, we will represent these operations as follows: For ciphertexts $\text{Enc}(m_1)$ and $\text{Enc}(m_2)$, which represent the encrypted plaintext messages m_1 and m_2 , respectively,

- $\text{Enc}(m_1) \oplus \text{Enc}(m_2) = \text{Enc}(m_1 + m_2)$
- $\text{Enc}(m_1) \odot m_2 = m_1 \odot \text{Enc}(m_2) = \text{Enc}(m_1 \cdot m_2)$
- $\text{Rot}(\text{Enc}(m_1); r) = \text{Enc}(\text{PRot}(m_1; r))$,

where $m_1 \cdot m_2$ denotes component-wise multiplication and $\text{PRot}(m_1; r)$ denotes the cyclically shifted plaintext vector of m_1 by r to the left. The RNS-CKKS is known to satisfy indistinguishability under chosen plaintext attack (IND-CPA) security [29], which guarantees that no information leakage occurs even after applying all supported homomorphic operations to ciphertexts. Moreover, it is well known that the RNS-CKKS homomorphic encryption scheme can also satisfy the stronger security notion of IND-CPA^D, simply by adding an appropriately sized Gaussian noise to the resulting ciphertexts while keeping the original homomorphic operations unchanged [30]. Since all algorithms proposed in this paper rely solely on the three main operations inherently supported by the RNS-CKKS scheme, the security of our algorithms is naturally ensured by the IND-CPA (or IND-CPA^D) security of RNS-CKKS. Moreover, the data integrity is preserved with only negligible error introduced by the approximate nature of the scheme, ensuring that the ciphertexts undergo minimal distortion during homomorphic computations.

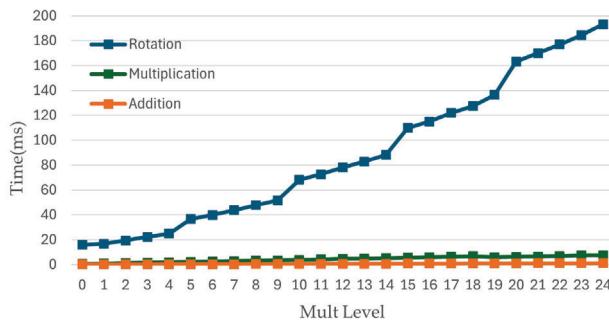


Figure 1: Latency of RNS-CKKS basic operation – addition, plaintext multiplication, rotation.

In this paper, the rotation operation of rotating a ciphertext by r is denoted as *rotation with r shift*. Rotation with the same shift requires the same rotation key during the operation. This means that rotation operations with different shifts require different rotation keys, and since the size of these rotation keys can be quite large, reducing them is one of the important challenges for FHE-based PPML. Also, it should be noted that in this paper, we only use plaintext multiplication, means that multiplication occurs only between plaintext and ciphertext, not between ciphertexts.

Every ciphertext has a unique positive integer, multiplicative level (abbreviated as mult level or level when there is no confusion). The level of a ciphertext indicates the number of remaining homomorphic multiplications that can be performed. When multiplying ciphertext with levels l_1 and l_2 , the resulting ciphertext will have a level of $\min(l_1, l_2) - 1$. Similarly, multiplication between l_1 level-ciphertext with plaintext, the resulting ciphertext will have a level of $l_1 - 1$. If the level becomes 0 due to repeated multiplication operations, further multiplication becomes impossible, and a bootstrapping operation is required to increase the level again. This operation is the most computationally intensive operation in RNS-CKKS. The number of levels consumed when performing an HE algorithm is referred to as the depth of the algorithm. Therefore, when designing a homomorphic encryption algorithm, it is crucial to consider not only the number of heavy operations like rotations and ciphertext-ciphertext multiplications, but also the depth consumed.

Latency of all homomorphic operations, including addition, multiplication, and rotation, are influenced by the level of the ciphertexts. Operations performed at higher levels require more time. The precise execution time of each operation varies slightly depending on several initial settings of CKKS, and in this paper, the time required for each operation in our environment is depicted in Figure 1.

As indicated in Figure 1, rotation operation takes the longest time compared to addition and plaintext multiplication operations, and it also exhibits the highest rate of increase in execution time with respect to the level. Therefore, to achieve optimization of execution time in programs

utilizing RNS-CKKS, it is effective to reduce the number of rotation operations.

3.2. Multiplexed Parallel Convolution

Due to its feature of supporting real number operations in an encrypted state, RNS-CKKS is widely used in Privacy-Preserving Machine Learning (PPML). In [8], CNN, particularly ResNet, was implemented using multiplexed parallel packing. Multiplexed parallel packing proposed to address situations involving convolutions with stride values greater than 1. It possesses the property of being multiplexed, meaning multiple channels are interleaved, and also exhibits the parallel property where the same data is replicated multiple times within a ciphertext.

The convolution algorithm, performed on ciphertext that is multiplexed parallel packed, is proposed as multiplexed parallel convolution (MPConv) in [8]. MPConv involves three main processes, denoted as SISOConv, RotationSum, and ZeroOutCombine. To express each process of MPConv, various parameters will be utilized, and precise definitions of each parameter are documented in Appendix B. SISOConv process is based on single-input single-output convolution proposed in [23, 31]. From the perspective of the size of plaintext, assume that the three-dimensional single input to this process is denoted as $A \in \mathbb{R}^{w_i \times h_i \times c_i}$. For single kernel K , which can be denoted as $K \in \mathbb{R}^{f_h \times f_w \times c_i}$, the result of SISOConv process is $B \in \mathbb{R}^{w_o \times h_o \times c_i}$. The RotationSum process corresponds to combining the values at each position of the channels in B . In this process, rotation and addition operations are repeated $\log_2 c_i$ times. The result of RotationSum process can be denoted as $C \in \mathbb{R}^{w_o \times h_o}$. In ZeroOutCombine process, since there are c_o numbers of kernels, combining c_o numbers of $C \in \mathbb{R}^{w_o \times h_o}$, we can get final results as $C' \in \mathbb{R}^{w_o \times h_o \times c_o}$. In particular, the ZeroOutCombine process can be viewed as divided into two subprocesses: the ZeroOut subprocess, which removes invalid values, and the Combine subprocess, which combines the results of each channel. Finally, to maintain the characteristic of the parallel existence of identical data within the ciphertext, the output is copied by a few rotation and addition operations.

Since rotation operations are relatively computationally intensive, the approximate performance of the algorithm can be gauged by the number of rotation operations performed. The number of rotations performed in each of the processes SISOConv, RotationSum, ZeroOutCombine and copying data are denoted as $f_h f_w - 1$, $q(2\lceil \log_2 k_i \rceil + \lceil \log_2 t_i \rceil)$, c_o , and $\log_2 p_o$, respectively. Considering that the number of rotations varies per level, let the execution time of rotation operations at level l be r_l . When multiplexed parallel convolution is executed on a ciphertext with level l' , and since the entire process consumes two multiplicative levels, the execution time can be summarized as follows. $r_{l'}(f_h f_w - 1) + r_{l'-1}(q(2\lceil \log_2 k_i \rceil + \lceil \log_2 t_i \rceil) + c_o) + r_{l'-2}(\log_2 p_o)$.

Throughout this paper, we use the notation in Figure 2 to illustrate a multiplexed parallel packed ciphertext, specifically for the case where the gap $k = 2$. Figure 3 simplifies the multiplexed parallel convolution where each variable is

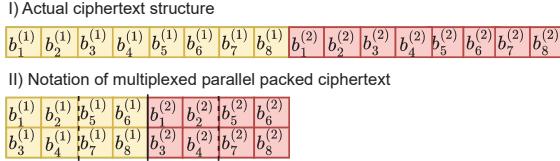


Figure 2: Notation of multiplexed parallel packed ciphertext in this paper.

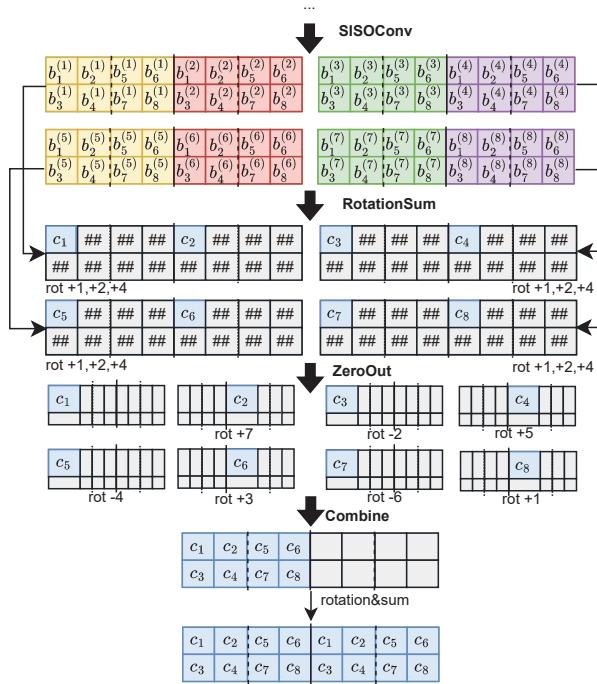


Figure 3: A conceptual diagram example of MPConv [8] where $w_i = h_i = w_o = h_o = 1$, $f_w = f_h = 3$, $c_i = c_o = 8$, $p_i = p_o = 2$, $k_i = k_o = 2$.

defined as $w_i = h_i = w_o = h_o = 1$, $f_w = f_h = 3$, $c_i = c_o = 8$, $p_i = p_o = 2$, $k_i = k_o = 2$. $b_j^{(i_1)}$ corresponds to the data of the j th channel when input data of convolution is multiplied by the weight of i_1 th kernel, and c_{i_1} can be defined as $c_{i_1} = \sum_{j=1}^8 b_j^{(i_1)}$.

In the input of the convolution, each ciphertext contains two identical parallelly existing data, and each data is multiplied by different kernel weights. This corresponds to the SISOConv process. In RotationSum process, each data has eight channels, so with $\log_2 8 = 3$ rotations per ciphertext, the values of the channels are summed. During this process, meaningless values are generated which are denoted as $\#\#$ in Figure 3. In the ZeroOutCombine process, these meaningless values are eliminated through multiplication operations and rearranged to the correct positions through rotation operations. Lastly, one rotation ensures that each

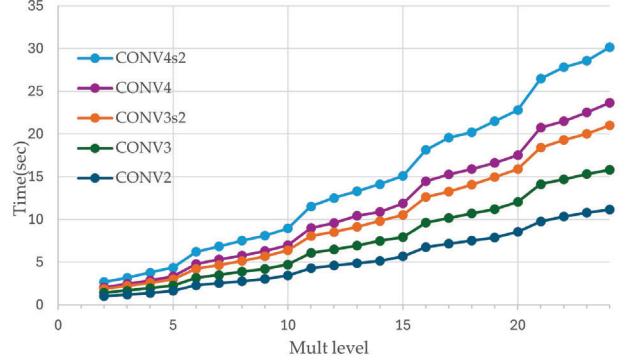


Figure 4: Latency of MPConv in various multiplication levels.

ciphertext maintains the structure of containing two identical parallelly existing data.

Furthermore, in conventional [8], multiplexed parallel convolution was tested in the situation of six convolutions, which are CONV1, CONV2, CONV3s2, CONV3, CONV4s2, and CONV4. To compare with this, our paper also conducted several tests targeting these convolutions. The parameters for each convolution can also be found in Appendix B. In [8], function Vec often used to map some three-dimension tensor to a vector to describe convolution. To avoid confusion, in our paper, Vec similarly defined as, for some three-dimension tensor $B \in \mathbb{R}^{h_i \times w_i \times c_i}$, $\text{Vec}(B)$ maps tensor B to a vector in \mathbb{R}^{n_t} . $\text{Vec}(B) = (b_0, \dots, b_{n_t-1}) \in \mathbb{R}^{n_t}$, where b can be defined as

$$b[i] = \begin{cases} B[\lfloor (i \bmod h_i w_i) / w_i \rfloor] & \\ [i \bmod w_i] [\lfloor i/h_i w_i \rfloor], & \text{if } 0 \leq i < h_i w_i c_i \\ 0, & \text{otherwise,} \end{cases}$$

4. Rotation-Optimized Multiplexed Parallel Convolution

As established in the Section 2, MPConv provides a broadly applicable and structurally flexible approach to homomorphic convolution. However, its practical deployment is hindered by significant computational bottlenecks, primarily due to its reliance on a large number of rotation operations. As illustrated in Figure 4, the execution time of MPConv grows super-linearly with the multiplicative level. Furthermore, the required rotation keys introduce a prohibitively large volume of data, creating a practical bottleneck for transmission and storage in client-server settings.

To overcome these fundamental limitations, this paper introduces Rotation-Optimized Multiplexed Parallel Convolution (RO-MPConv). The core principle of our method is to restructure the convolution algorithm to minimize the number of rotation operations, thereby reducing both computational latency and rotation key volume without sacrificing other performance metrics. This section details the three key strategies behind RO-MPConv: a fundamental redesign

to reduce rotations, the exploitation of the correlation between the RotationSum and CrossCombine processes, and the strategic use of additional multiplicative depth for further performance trade-offs.

4.1. High-Level Idea

The primary strategy of RO-MPConv is to reduce both computational latency and key transmission volume by minimizing the number of rotation operations. As summarized in Table 1, the ZeroOutCombine process in the conventional MPConv is a prime target for optimization. This is because the number of rotation operations and the number of required rotation keys are nearly identical, implying that each rotation uses a unique shift and, therefore, a unique rotation key. Consequently, reducing the rotations in this specific process yields a dual benefit: decreased latency and a smaller key volume.

Table 1
Number of rotation keys and rotation operations during MPConv.

Convolution	Process	# RotationKey	# Rotation operation
CONV1	SISOConv	8	8
	RotationSum	2	4
	ZeroOutCombine	16	17
CONV2	SISOConv	8	8
	RotationSum	4	32
	ZeroOutCombine	16	17
CONV3s2	SISOConv	8	8
	RotationSum	4	64
	ZeroOutCombine	33	34
CONV3	SISOConv	8	8
	RotationSum	5	40
	ZeroOutCombine	33	34
CONV4s2	SISOConv	8	8
	RotationSum	5	80
	ZeroOutCombine	66	67
CONV4	SISOConv	8	8
	RotationSum	6	48
	ZeroOutCombine	66	67

The key idea of our approach is to modify the order of convolution weights and the direction of rotations during RotationSum process by considering the final arrangement of the multiplexed parallel packed result from the outset. In the conventional MPConv, illustrated in Figure 3, the data corresponding to each channel is aligned to the same position (the top-left in the diagram) after the RotationSum process. However, to restore the multiplexed parallel packed state in the final output, the Combine subprocess must then perform a number of rotation operations nearly equal to the number of output channels to move each channel to its correct final location.

Figure 5 illustrates the high-level idea of RO-MPConv, where the entire process is conducted in consideration of

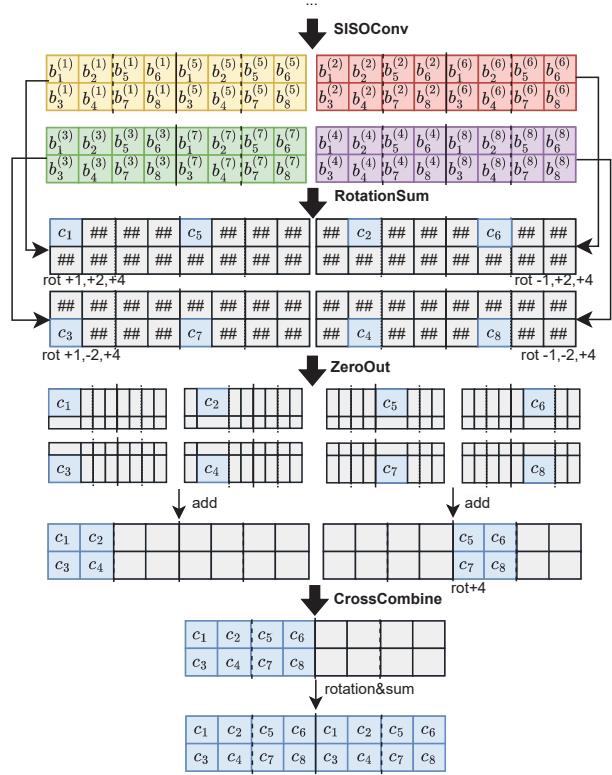


Figure 5: A conceptual diagram of the high-level idea of RO-MPConv, under the same conditions as in Figure 3.

the final channel arrangement. In the final output ciphertext, channels c_1 through c_4 are grouped in the initial block of the ciphertext, while channels c_5 through c_8 are grouped in the subsequent block. Furthermore, within these blocks, channel pairs such as (c_1, c_5) , (c_2, c_6) , (c_3, c_7) , (c_4, c_8) occupy the same relative positions.

To achieve this structure efficiently, we introduce two key modifications. First, the order in which convolution weights are multiplied in the SISOConv process is revised. Instead of a simple sequential multiplication, weights are applied in an order that pre-positions the resulting channels. For example, the results for channels 1-4 are generated from the first parallel data block, while channels 5-8 are generated from the second. Additionally, channel pairs with similar target positions in the final layout (e.g., 1 and 5, 2 and 6) are generated within the same ciphertext. This deliberate placement ensures that they can be processed with the same rotation direction in the subsequent RotationSum stage. Second, during the RotationSum process, ciphertexts undergo rotations in both positive and negative directions, not just positive ones. This strategic variation in rotation direction allows the resulting valid values (the summed channels) to be placed in locations that closely approximate their final positions in the output ciphertext, rather than being gathered in a single default location.

As a result of these modifications, the majority of channels are already aggregated into near-final positions before

the final combining step. This allows us to replace the original Combine process, which meticulously repositioned each channel one by one, with a new, far more efficient CrossCombine subprocess. The CrossCombine operation takes these large, pre-arranged blocks of channels and moves them into their final configuration with only a minimal number of rotations.

It is important to note that while the modifications in the SISOConv and RotationSum processes are crucial for enabling this optimization, they are performance-neutral. The actual performance gain arises exclusively from replacing the Combine process with the CrossCombine process. Quantitatively, this change reduces the required number of rotation operations from c_o to $p_c - 1$. In the six convolution scenarios previously mentioned, c_o is approximately 16, 32, or 64, whereas p_c is much smaller value, around 2, 4, or 8. This leads to a notable observation: RO-MPConv is particularly effective for convolutions with a large number of output channels (c_o) as the reduction in rotations becomes more significant as (c_o) increases.

This improvement is reflected in the revised execution time formulation. Defining the execution time of a rotation at level l as r_l , the total execution time for RO-MPConv operating on a ciphertext at level l' , the execution time performance can be summarized as follows: $r_{l'}(f_h f_w - 1) + r_{l'-1}(q(\lceil \log_2 k_i \rceil + \lceil \log_2 t_i \rceil)) + r_{l'-2}(p_c - 1 + \log_2 p_o)$.

4.2. Using correlation of RotationSum process and CrossCombine process

The objectives of the three main processes in multiplexed parallel convolution are distinct. The SISOConv process computes the product of the input data and the kernel weights. The RotationSum process then aggregates the values from each of these resulting channels. Finally, the ZeroOut process filters out invalid values generated during RotationSum and rearranges the valid values into their correct final positions.

As established in the previous section, the replacement of the Combine subprocess with CrossCombine in RO-MPConv reduces the number of rotations required in the final stage. This newfound efficiency opens up an opportunity for a further optimization: a strategic trade-off between the RotationSum and CrossCombine processes. The core idea is to perform fewer merge operations during the RotationSum phase and compensate for the incomplete summation by performing additional rotations during the CrossCombine phase.

In the conventional RotationSum process, each rotation-and-addition operation typically halves the number of distinct channel components within a ciphertext, requiring $\lceil \log_2 c_i \rceil$ rotations per ciphertext to fully merge all channels. If we reduce the number of rotations performed per ciphertext by x , the number of rotations in the RotationSum stage decreases by $q \cdot x$. However, this leaves 2^x partial sums within each ciphertext that must be combined later. Consequently, the number of parallel data blocks to be merged in the

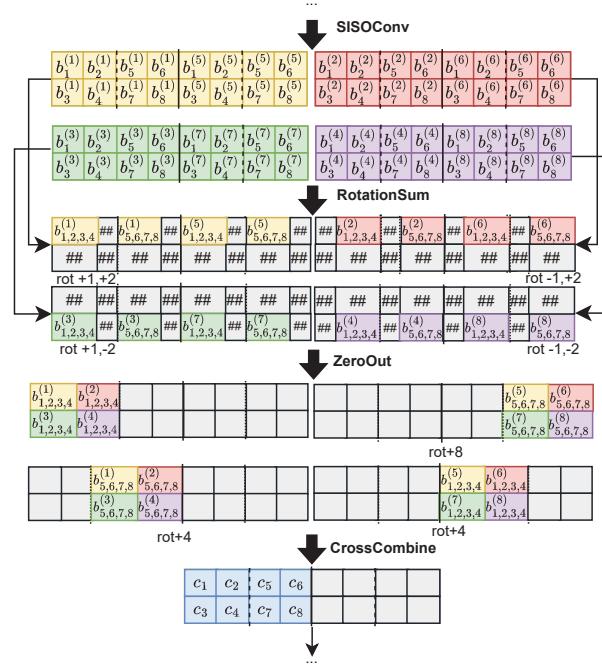


Figure 6: A conceptual diagram showing how RO-MPConv leverages the correlation between the RotationSum and CrossCombine stages to reduce rotation counts.

CrossCombine stage increases by a factor of 2^x . This trade-off is generally infeasible in the original MPConv because the number of rotations in its Combine step is already large (c_o), and doubling it would negate any savings. In contrast, RO-MPConv's CrossCombine process requires only a small number of rotations ($p_c - 1$), making this trade-off viable and often advantageous.

Figure 6 provides a conceptual example of this technique. In this diagram, the notation $b_{j_1,j_2,j_3,j_4}^{(i)}$ represents the sum of four channel values: $b_{j_1}^{(i)} + b_{j_2}^{(i)} + b_{j_3}^{(i)} + b_{j_4}^{(i)}$. After the SISOConv process, each parallel data block contains 8 channels. In the standard RO-MPConv shown in Figure 5, the RotationSum process would perform $3 \times 4 = 12$ rotations, followed by 1 rotation in the CrossCombine process, for a total of 13 rotations. However, in Figure 6, we reduce the number of rotations in RotationSum to 2 per ciphertext, for a total of $2 \times 4 = 8$ rotations. This leaves two partial sums per block. The number of items to be merged in CrossCombine, therefore, increases from 2 to 4, requiring 3 rotations. The new total is $8 + 3 = 11$ rotations, demonstrating a net reduction.

However, the optimal number of rotations to defer is not fixed, as the execution time of a rotation varies with its multiplicative level. Specifically, rotations in the RotationSum process occur at a higher (and thus slower) level than those in the CrossCombine process. Therefore, a careful balance must be struck to achieve the minimum possible latency.

To formalize this trade-off, let x denote the number of rotation stages deferred from the RotationSum process. When using this correlation, the optimal value of x for a 2-depth convolution starting at level l can be found by solving the following minimization problem:

$$\operatorname{argmin}_x (r_{l-1} q (\lceil \log_2 c_i \rceil - x) + r_{l-2} (2^x p_c - 1)) \quad (1)$$

where $x \in [0, \log_2 (2^{\lceil \log_2 c_i \rceil} p_i / c_o)]$. The first term in the equation accounts for the total time of the reduced rotations in RotationSum, while the second term reflects the time for the increased rotations in CrossCombine. The upper bound for x is determined by the finite length of the ciphertext, which limits the number of partial sums that can be combined, an issue discussed further in Section 8.

4.3. Convolution with additional depth

The latency of FHE algorithms is predominantly influenced by the number of computationally heavy homomorphic operations, such as rotations and ciphertext multiplications, as well as the frequency of bootstrapping required when the multiplicative level is exhausted. Therefore, designing efficient FHE algorithms necessitates a careful consideration of the trade-off between the count of these heavy operations and the consumed multiplicative depth. This design principle is a well-established paradigm in FHE literature, where variants of operations are proposed to navigate this balance. For instance, [32] organizes the trade-off of consuming additional depth and execution time in SlotToCoeff and CoeffToSlot operations in bootstrapping, and [33] organizes the trade-off of consuming additional depth and accuracy in relu operation.

In this section, we introduce a family of RO-MPConv variants designed to trade additional multiplicative depth for further latency reduction. This flexibility, which is not available in the original MPConv, offers a new performance trade-off. The key idea of this technique is to combine ciphertexts that require the same rotation during the RotationSum process. When looking at the conventional multiplexed parallel convolution in Figure 3, this kind of optimization is challenging to employ. Because all ciphertexts undergo the same rotation operation during the RotationSum process, their valid values occupy identical positions, making it difficult to combine them without additional operations.

In contrast, in Figure 5, rotation-optimized convolution rotates ciphertexts in different directions during the RotationSum processes, resulting in varied positions where valid values exist. Therefore, leveraging this aspect, optimization as Figure 7 can be achieved. In the RotationSum process, ciphertexts requiring a rotation of +2 and those needing -2 are combined separately and then subjected to rotation operations. Since the valid positions of each ciphertext are different, no additional operations are necessary. However, to combine each ciphertext, invalid values must be filtered out, necessitating the use of multiplication operations and an additional level.

For an exact comparison, compared to the 2-depth consuming RO-MPConv in Figure 5, the number of rotations

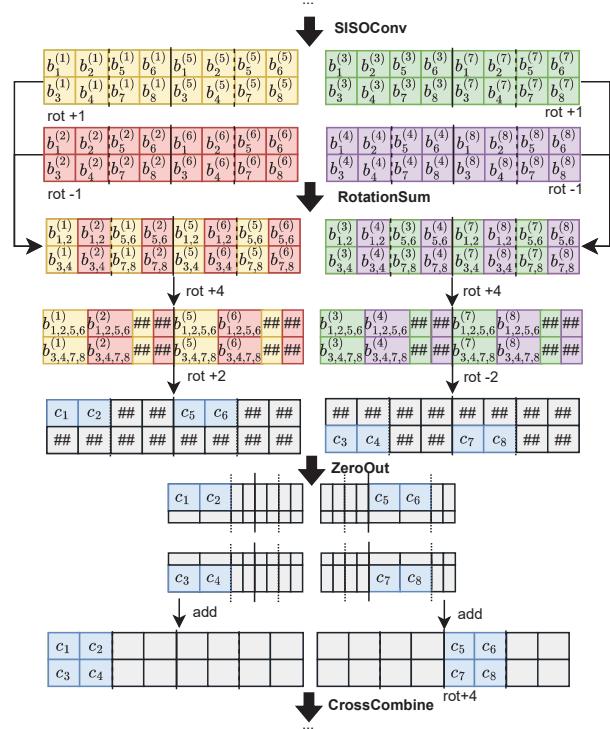


Figure 7: Conceptual diagram of the RO-MPConv variant that utilizes an additional multiplicative depth.

used during RotationSum process and ZeroOut process reduced by 13 to 9. It can be observed that using additional depth can reduce the number of rotation operations. Moreover, from the perspective of rotation keys, when combining ciphertexts that require the same shift and then performing rotation operations, the number of rotation operations decreases without altering the shifts of rotation operations. Consequently, using a RO-MPConv with additional depth has a negligible impact on the total number of rotation keys.

The maximum additional depth that an RO-MPConv variant can consume is a function of q , the number of ciphertexts available for merging in the RotationSum process. The underlying mechanism involves iteratively combining these q ciphertexts in pairs, where each iteration of merging consumes one additional multiplicative level. This process can be repeated a maximum of $\log_2 q - 1$ times before the ciphertexts are fully consolidated. Given that the standard RO-MPConv already consumes two levels, the maximum total depth for a variant becomes $2 + (\log_2 q - 1) = \log_2 q + 1$

Similar to the approach in Subsection 4.1, we can now formalize the performance of RO-MPConv variants that utilize both the RotationSum-CrossCombine correlation and the additional depth-consuming technique. If a convolution consuming $d+1$ ($d \geq 2$) depth is utilized, starts at level l ($l \geq d+1$) then the approximate performance of this convolution in the RotationSum, and CrossCombine processes can be

expressed as the following equation:

$$r_{l-1}(x_1 q) + \sum_{m=2}^d r_{l-m} \left(\frac{x_m q}{\prod_{j=1}^{m-1} q_j} \right) + r_{l-d-1} (2^x p_c - 1) \quad (2)$$

where variables defined as $\prod_{m=1}^d q_m = q$, $\lceil \log_2 c_i \rceil - x = \sum_{m=1}^d x_m$ and $x_m \geq \log_2 q_m$ for all m . x represents the same as in Subsection 4.2. In other words, $\lceil \log_2 c_i \rceil - x$ denotes the number of rotation operations needed for RotationSum process. In the RotationSum process, the variable determining how many rotations should be performed before the m -th combining of the ciphertexts is denoted as x_m , and the number of ciphertexts being combined is denoted as q_m . Since x_m rotation and addition combines 2^{x_m} data in ciphertext, maximum 2^{x_m} ciphertexts can be combined during m -th combining. Therefore, value of q_m is bounded by x_m as $x_m \geq \log_2 q_m$.

Providing guidelines to use additional depth in RO-MPConv, first, the value of d , which specifies how many depths will be consumed during convolution, has to be determined. Then, while satisfying the following conditions $\prod_{m=1}^d q_m = q$, $\lceil \log_2 c_i \rceil - x = \sum_{m=1}^d x_m$ and $x_m \geq \log_2 q_m$ for all m , x , x_m and q_m within the range $1 \leq m \leq d$ to minimize Equation 2 has to be selected. Since the number of combinations satisfying all conditions and the simplicity of Equation 2 allow for brute force, finding its minimum value is not a burdensome task. Moreover, x , x_m and q_m are values pre-determined during the implementation process of RO-MPConv, not during the operation itself, meaning the process of finding the minimum value of Equation 2 has no impact on runtime performance.

4.4. Algorithm Description

In presenting the main idea of RO-MPConv, we focused on reducing the number of rotations required in the ZeroOutCombine process to reduce execution time and the size of rotation keys. Also, by leveraging the relationship between RotationSum and CrossCombine, as well as using convolution with additional depth, we were able to achieve further execution time reduction. In this section, the focus is on explaining the algorithm of RO-MPConv, taking into account such ideas.

Through our high-level idea proposed in Subsection 4.1, it is evident that optimization of computations is feasible. However, due to the nature of the idea, which involves considering the arrangement of results, the method of carrying out computations is not unique. For instance, Figure 8 is a modified convolution of Figure 5. Upon observation, it can be seen that although the batch of kernels and types of rotations may vary slightly, the output of the convolution remains the same, and the number of rotations is also identical. These modifications in the RO-MPConv algorithm do not affect the performance, which means the number of rotations or size of rotation keys doesn't change and it only changes the shift of rotation operation. Indeed, the RO-MPConv algorithm is not unique and possesses a heuristic aspect.

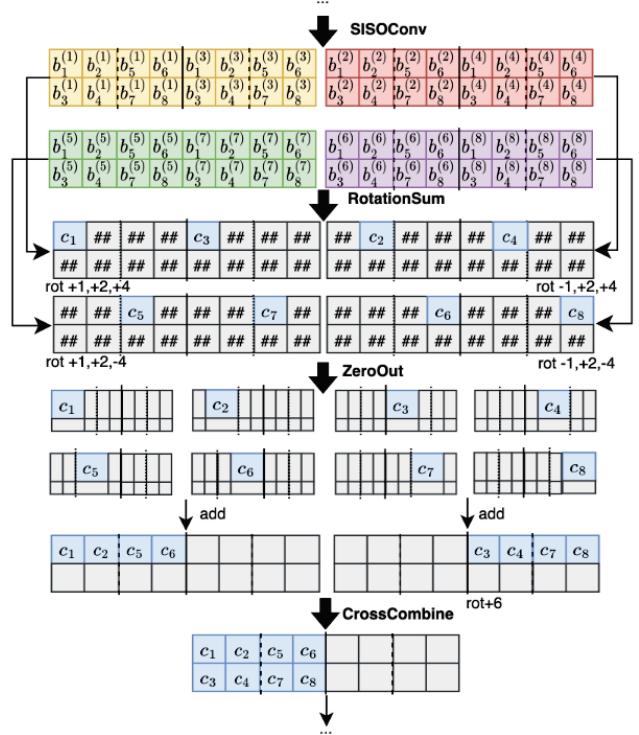


Figure 8: Showing the non-uniqueness property of RO-MPConv. Same condition, and same performance as Figure 5, but the kernel arrangement and shift amount of rotation are different.

Thus, we will not statically generalize how kernel weights are arranged and what rotation operations are employed in the RotationSum process and CrossCombine process; rather, we will provide **blueprints** in detailing these aspects. Each rotation-optimized convolution requires three blueprints: KernelBP, RotationSumBP, and CrossCombineBP. These blueprints contain information about the kernel placement, the rotation shifts used during the RotationSum process, and the configuration for the CrossCombine process, respectively. The precise values for each blueprint used in our experiments can be found in Appendix A.

To describe the proposed RO-MPConv algorithm, we first define several utility functions and tensors for manipulating data within ciphertexts. These include valid value-selecting tensors (S, S', S'') and a weight-to-vector mapping function, (ParWgt).

$S^{(r)}$ filter serves to select a valid value from the ciphertext, which is rotated by r . $S^{(r)} = (S^{(r)}[i_1])_{0 \leq i_1 < n_t} \in \mathbb{R}^{n_t}$. The components of filter $S^{(r)}$ are defined as follows:

$$S^{(r)}[i_1] = \begin{cases} 1, & \text{if } (r > 0 \text{ and } (i_1 \bmod 2r) < r) \text{ or} \\ & (r < 0 \text{ and } (i_1 \bmod |2r|) \geq |r|) \\ 0, & \text{otherwise,} \end{cases} \quad (3)$$

$S'^{(k,l)}$ filters out parallel data in the ciphertext. For ciphertext which has l identical data parallelly, filter $S'^{(k,l)}$ select k -th data. $S'^{(k,l)} = (S'^{(k,l)}[i_1])_{0 \leq i_1 < n_t} \in \mathbb{R}^{n_t}$. The components of

filter $S'^{(k,l)}$ are defined as follows:

$$S'^{(k,l)}[i_1] = \begin{cases} 1, & \text{if } n_t k/l \leq i_1 < n_t(k+1)/l \\ 0, & \text{otherwise,} \end{cases} \quad (4)$$

S'' filters out data considering when the stride value of convolution is bigger than 1. $S'' = (S''[i_1])_{0 \leq i_1 < n_t} \in \mathbb{R}^{n_t}$. The components of the filter S'' are defined as follows:

$$S''[i_1] = \begin{cases} 1, & \text{if } i_1 \bmod k_o^2 w_o < k_i k_o w_o \text{ and } i_1 \bmod k_o < k_i \\ 0, & \text{otherwise,} \end{cases} \quad (5)$$

The function $\text{ParWgt}(U; i_1, i_2, i_3)$ maps the weight tensor $U \in \mathbb{R}^{h_i \times w_i \times c_i \times c_o}$ to a plaintext vector in \mathbb{R}^{n_t} via $\text{ParWgt}(U; i_1, i_2, i_3) = \text{Vec}(\overline{U}''^{(i_1, i_2, i_3)})$ where $i_1 \in [0, f_h]$, $i_2 \in [0, f_w]$, $i_3 \in [0, q]$. Parallelly multiplexed shifted weight tensor $\overline{U}''^{(i_1, i_2, i_3)}$ is defined as follow. $\overline{U}''^{(i_1, i_2, i_3)} = (\overline{U}''^{(i_1, i_2, i_3)}[i_5][i_6][i_7]) \in \mathbb{R}^{k_i h_i \times k_i w_i \times t_i p_i}$ where $i_5 \in [0, k_i h_i]$, $i_6 \in [0, k_i w_i]$, and $i_7 \in [0, t_i p_i]$.

$$\overline{U}''^{(i_1, i_2, i_3)}[i_5][i_6][i_7] = \begin{cases} 0, & \text{if } k_i^2(i_7 \bmod t_i) + k_i(i_5 \bmod k_i) \\ & + i_6 \bmod k_i \geq c_o \\ & \text{or } \lfloor i_7/t_i \rfloor + p_i i_3 \geq c_o \\ & \text{or } \lfloor i_5/k_i \rfloor - (f_h - 1)/2 \\ & + i_1 \notin [0, h_i - 1] \\ & \text{or } \lfloor i_6/k_i \rfloor - (f_w - 1)/2 \\ & + i_2 \notin [0, w_i - 1], \\ U[i_1][i_2][i_8][i_9] & \text{otherwise,} \end{cases} \quad (6)$$

i_8 and i_9 are defined as $i_8 = k_i^2(i_7 \bmod t_i) + k_i(i_5 \bmod k_i) + (i_6 \bmod k_i)$, and $i_9 = \text{KernelBP}[p_i i_3][\lfloor i_7/t_i \rfloor]$.

Since RotationSum process used in RO-MPConv has different rules, we redefine the SumSlots algorithm in [8] to AdaptSumSlots, which contains RotationSum process and also SISOConv.

The execution flow of Algorithm 1 is structured according to the provided blueprints, which define operational details such as the opType that specifies the type of operation. If opType is 0, for example, the blueprint might be [0, 2048], indicating rotating the ciphertext by 2048 and then performing an addition operation. When opType is 1, it signifies performing a Combine operation in the middle of RotationSum, meaning using additional depth. An example blueprint for this would be [1, 4, 1, 32], indicating rotating each of the four ciphertexts by +1+32, -1+32, +1-32, and -1-32, respectively, then filtering out invalid values and combining them. When opType is 2, it represents operations before the CrossCombine stage after the RotationSum process. Similar to opType 1, in this case, multiplication is carried out as many times as the number of data involved in the CrossCombine during the filtering process to classify

Algorithm 1 AdaptSumSlots(i' , d , ct'_α , U)

Input: Indicating the order of ciphertext i' , Rotated parallelly multiplexed tensor ciphertexts ct'_α , Current location within RotationSumBP d , and weight tensor U

Output: Tensor ciphertext ct_b and filter of ct_b that apply zero out S_b .

```

1:  $\text{ct}_b \leftarrow \text{ct}_{\text{zero}}$ 
2:  $S_b \leftarrow S_{\text{one}}$ 
3:  $\text{opType} \leftarrow \text{RotationSumBP}[d][0]$ 
4: if  $d = 0$  then
5:   for  $i_1 \leftarrow 0$  to  $f_h - 1$  do
6:     for  $i_2 \leftarrow 0$  to  $f_w - 1$  do
7:        $\text{ct}_b \leftarrow \text{ct}_b \oplus \text{ct}'_\alpha^{(i_1, i_2)} \odot \text{ParWgt}(U; i_1, i_2, i')$ 
8:     end for
9:   end for
10: end if
11: if  $d > 0$  then
12:   if  $\text{opType} = 0$  then
13:      $\text{ct}_a, S_a \leftarrow \text{AdaptSumSlots}(i', d - 1, \text{ct}'_\alpha, U)$ 
14:      $r' \leftarrow \text{RotationSumBP}[d][1]$ 
15:      $\text{ct}_b \leftarrow \text{ct}_a \oplus \text{Rot}(\text{ct}_a; r')$ 
16:      $S_b \leftarrow S_a \odot S^{(r')}$ 
17:   end if
18:   if  $\text{opType} = 1$  then
19:      $\text{SumN} \leftarrow \text{RotationSumBP}[d][1]$ 
20:     for  $i'' \leftarrow \text{SumN} \times i'$  to  $(\text{SumN} + 1) \times i' - 1$  do
21:        $\text{ct}_a, S_a \leftarrow \text{AdaptSumSlots}(i'', d - 1, \text{ct}'_\alpha, U)$ 
22:       for  $j \leftarrow 2$  to  $\log_2(\text{SumN}) + 1$  do
23:          $r' \leftarrow \text{RotationSumBP}[d][j]$ 
24:         if  $((i'' >> (j - 2)) \& 1) = 1$  then
25:            $\text{ct}_a \leftarrow \text{ct}_a \oplus \text{Rot}(\text{ct}_a; -r')$ 
26:            $S_a \leftarrow S_a \odot S^{(-r')}$ 
27:         else
28:            $\text{ct}_a \leftarrow \text{ct}_a \oplus \text{Rot}(\text{ct}_a; r')$ 
29:            $S_a \leftarrow S_a \odot S^{(r')}$ 
30:         end if
31:       end for
32:        $\text{ct}_b \leftarrow \text{ct}_b \oplus (\text{ct}_a \odot S_a)$ 
33:     end for
34:   end if
35:   if  $\text{opType} = 2$  then
36:      $\text{ct}_b, S_b \leftarrow \text{AdaptSumSlots}(i', d - 1, \text{ct}'_\alpha, U)$ 
37:     for  $j \leftarrow 2$  to  $\log_2(\text{RotationSumBP}[d][1]) + 1$  do
38:        $r' \leftarrow \text{RotationSumBP}[d][j]$ 
39:       if  $((i' >> (j - 2)) \& 1) = 1$  then
40:          $\text{ct}_b \leftarrow \text{ct}_c \oplus \text{Rot}(\text{ct}_b; -r')$ 
41:          $S_b \leftarrow S_b \odot S^{(-r')}$ 
42:       else
43:          $\text{ct}_b \leftarrow \text{ct}_c \oplus \text{Rot}(\text{ct}_b; r')$ 
44:          $S_b \leftarrow S_b \odot S^{(r')}$ 
45:       end if
46:     end for
47:      $S_b \leftarrow S_b \odot S''$ 
48:   end if
49: end if
50: return  $\text{ct}_b, S_b$ 

```

Algorithm 2 RO-MPConv(ct'_a, U)

```

1: Input: Parallelly multiplexed tensor ciphertext  $\text{ct}'_a$  and
   weight tensor  $U$ .
2: Output: Parallelly multiplexed tensor ciphertext  $\text{ct}'_c$ 
3:  $\text{ct}'_c \leftarrow \text{ct}_{\text{zero}}$ 
4: for  $i_1 \leftarrow 0$  to  $f_h - 1$  do
5:   for  $i_2 \leftarrow 0$  to  $f_w - 1$  do
6:      $\text{ct}'_{a(i_1,i_2)} \leftarrow \text{Rot}(\text{ct}'_a; k_i^2 w_i (i_1 - (f_h - 1)/2) + k_i (i_2 - (f_w - 1)/2))$ 
7:   end for
8: end for
9:  $\text{CrossNum} \leftarrow \text{CrossCombineBP}[0]$ 
10: for  $i_3 \leftarrow 0$  to  $\text{CrossNum} - 1$  do
11:    $\text{ct}'_{\beta(i_3)} \leftarrow \text{ct}_{\text{zero}}$ 
12: end for
13:  $\text{BPlen} \leftarrow \text{RotationSumBP}[0][0]$ 
14: for  $i_4 \leftarrow 0$  to  $\text{RotationSumBP}[\text{BPlen}][1] - 1$  do
15:    $\text{ct}'_b, S_a \leftarrow \text{AdaptSumSlots}(i_4, \text{BPlen}, \text{ct}'_a, U)$ 
16:   for  $i_3 \leftarrow 0$  to  $\text{CrossNum} - 1$  do
17:      $\text{ct}'_{\beta(i_3)} \leftarrow \text{ct}'_{\beta(i_3)} \oplus \text{ct}'_b \odot (S_a \odot S'^{(i_3, \text{CrossNum})})$ 
18:   end for
19: end for
20: for  $i_1 \leftarrow 0$  to  $\text{CrossNum} - 1$  do
21:    $\text{ct}'_c \leftarrow \text{ct}'_c \oplus \text{Rot}(\text{ct}'_{\beta(i_1)}; \text{CrossCombineBP}[i_1 + 1])$ 
22: end for
23: for  $j \leftarrow 0$  to  $\log_2 p_o - 1$  do
24:    $\text{ct}'_c \leftarrow \text{ct}'_c \oplus \text{Rot}(\text{ct}'_c; -2^j (n_t/p_o))$ 
25: end for
26: return  $\text{ct}'_c$ 
```

each data separately in parallel. Furthermore, as a detailed rule, $\text{RotationSumBP}[0][0]$ stores the maximum index of RotationSumBP , and $\text{RotationSumBP}[i][0]$ contains information about the i -th opType. If opType is 1 or 2, then $\text{RotationSumBP}[i][1]$ contains information about how many ciphertexts will be merged.

For RO-MPConv algorithm, the CrossCombineBP is used, indicating a simple form where each ciphertext at the i -th position is rotated by $\text{CrossCombineBP}[i]$ during CrossCombine process. As a detail, $\text{CrossCombineBP}[0]$ stores the number of data to be combined in CrossCombine process. For an intuitive example, consider the operation illustrated in Figure 8. In this case, the KernelBP is $[[1,3],[2,4],[5,7],[6,8]]$, which reflects the pairing of kernels within each ciphertext. The RotationSum process is then guided by its own blueprint, which specifies a common rotation with a shift of +2 for all ciphertexts, while rotations with shifts of ± 1 or ± 4 are applied variably to different ciphertexts. Moreover, after rotations with shifts of ± 1 or ± 4 , ZeroOut should be applied, hence opType becomes 2. Thus, $\text{RotationSumBP} = [[2],[0,4],[2,4,1,4]]$. In the CrossCombine process, since the second data is combined with a rotation operation of shift +6, $\text{CrossCombineBP} = [2,0,6]$.

A comparison of the performance of RO-MPConv and MPConv centered around the number of rotation operations can be summarized as follows. Still, the time required for rotation operations at mult level l' of the ciphertext is denoted as $r_{l'}$. For conventional MPConv, assuming that operations start at level l , its performance can be represented as follows: $r_l(f_h f_w - 1) + r_{l-1}(q(2\lceil \log_2 k_i \rceil + \lceil \log_2 t_i \rceil) + c_o) + r_{l-2}(\log_2 p_o)$.

For 2-depth consuming RO-MPConv which prioritizing minimizing execution time, assuming that operations start at level l , its performance can be represented as follows: $\min_x(r_l(f_h f_w - 1) + r_{l-1}(q(\lceil \log_2 c_i \rceil - x) + r_{l-2}(2^x p_c - 1 + \log_2 p_o))$ where $x \in [0, \log_2(2^{\lceil \log_2 c_i \rceil} p_i / c_o)]$.

For 3 or more depth consuming RO-MPConv which prioritize minimizing execution time, assuming that operations start at level l , its performance can be represented as follows: $\min_{x,x_1,\dots,x_d,q_1,\dots,q_m}(r_l(f_h f_w - 1) + r_{l-1}(x_1 q) + \sum_{m=2}^d r_{l-m} \left(x_m \frac{q}{\prod_{j=1}^{m-1} q_j} \right) + r_{l-d-1}(2^x p_c - 1 + \log_2 p_o))$ where $x \in [0, \log_2(2^{\lceil \log_2 c_i \rceil} p_i / c_o)]$, $\prod_{m=1}^d q_m = q$, $\lceil \log_2 c_i \rceil - x = \sum_{m=1}^d x_m$ and $x_m \geq \log_2 q_m$ for all m .

Once again, due to the heuristic nature of our approach, this paper does not provide a single, fixed algorithm for RO-MPConv. Instead, we have presented a set of performance equations that define the expected optimal performance for any given convolution configuration. These equations, therefore, serve as a benchmark for validation; an implementation of RO-MPConv can be considered correct if its performance aligns with the predictions derived from these equations.

5. Parallel BSGS Matrix-Vector Multiplication

Matrix-vector multiplication, including inference of fully connected layers of neural networks, is a common operation in recent AI models. Several methods have been devised for FHE matrix-vector multiplication [12, 34, 35, 36, 37]. One of the well-known methods to multiply a plaintext matrix and a ciphertext vector is the diagonal method from [12], and applying the Baby-step Giant-step(BSGS) algorithm to further reduce the number of rotations to reduce latency. [8] also implemented the fully connected layer by applying the BSGS algorithm to the diagonal method (hereinafter referred to as the BSGS diagonal method).

The characteristic of the BSGS diagonal method is that it determines the number of slots used only based on the lengths of the matrix and vector, regardless of the length of the ciphertext. For instance, when implementing ResNet targeting the CIFAR-10 ([38]) images in [8], the length of the ciphertext used is 32768. The fully connected layer of this ResNet involves a matrix-vector multiplication with a matrix of size 10×64 and a ciphertext containing 64 valid values. Due to the characteristics of multiplexed parallel packing, this operation starts with eight identical data within a single ciphertext. However, the BSGS diagonal method does not take this into account and operates only on one data, using

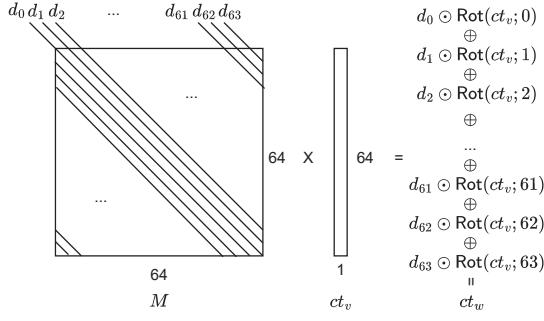


Figure 9: A simple illustration of diagonal method([12]) for matrix-vector multiplication.

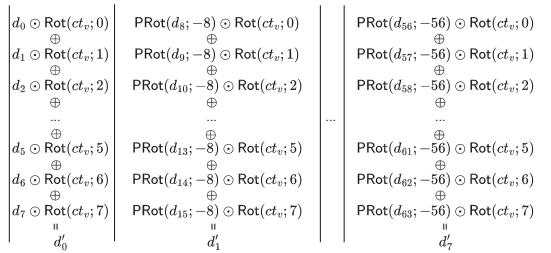


Figure 10: A simple illustration of BSGS diagonal method.

only around 72 slots for the matrix-vector multiplication. In other words, it does not efficiently utilize the ciphertext.

Since underutilizing the ciphertext usually necessitates additional computations, many FHE algorithms, including the conventional method [8], frequently employ structures with multiple identical data within a single ciphertext. This pattern is particularly evident in bootstrapping, the most computationally intensive operation in FHE. Sparse slot bootstrapping [39], widely adopted across many FHE systems, operates under the assumption that multiple identical data exist within a single ciphertext.

Therefore, in this section, we propose parallel BSGS matrix-vector multiplication, which supports the multiplication of a single plaintext square matrix and ciphertext vector containing multiple identical data. Parallel BSGS matrix-vector multiplication optimizes the number of rotations required compared to BSGS diagonal method, especially when ciphertext can contain numerous identical data parallel, including multiplexed parallel packed ciphertext. To explain parallel BSGS matrix-vector multiplication, understanding of both the conventional diagonal method and the BSGS diagonal method is required. For an intuitive explanation, we will offer images for an example of the multiplication of an $n \times n$ matrix and an $n \times 1$ ciphertext where $n = 64$. Here, the matrix is represented as M , the $n \times 1$ vector as v , and the result as w . In other words, $w = Mv$. Since w and v are ciphertexts in practice, we denote the ciphertext containing

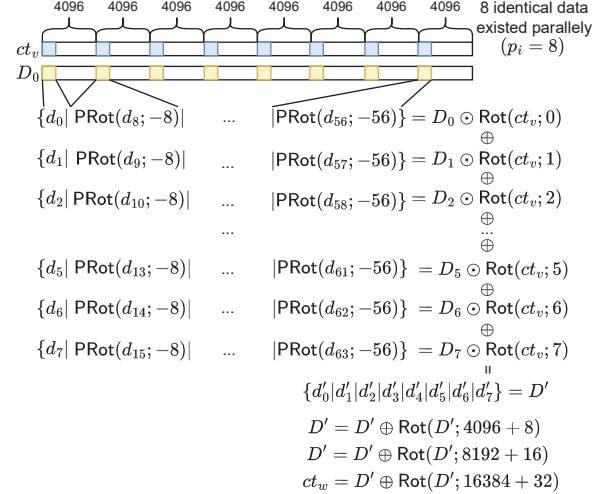


Figure 11: A simple illustration of parallel BSGS diagonal method.

the vector v at its head as ct_v , and the ciphertext for w as ct_w . Both ct_v and ct_w has the same length, n_t .

Figure 9 illustrates the diagonal method introduced in [12]. For diagonals of M , d is defined as $d_i[j] = A[j][j+i]$ where $i, j \in [0, n]$. And each d_i 's length is n_t , which is the same as ct_v . As shown in Figure 9, we can get ct_w by integrating the multiplication of each d_i and $\text{Rot}(ct_v; i)$. More precisely, it can be expressed as follows. $ct_w = \bigoplus_{i=0}^{n-1} (d_i \odot \text{Rot}(ct_v; i))$. Since the number of rotations and plaintexts needed during the diagonal method is determined by the number of d_i , it can be observed that $n - 1 = 63$ rotations and $n = 64$ plaintexts are required.

By applying the BSGS algorithm to Figure 9, we can examine examples of BSGS diagonal method as Figure 10. The idea leverages the fact that d_i 's are plaintexts, so pre-computing them incurs no overhead. Therefore, optimization is achieved by pre-rotating d_i 's. To represent BSGS diagonal method, several additional parameters are defined. n_1 denotes the small number of pre-rotated ciphertexts, where $n = n_1 n_2$. d' denotes the intermediate sums of products between pre-rotated plaintexts and ciphertexts. $d'_i = \bigoplus_{j=n_1 i}^{n_1 i + n_1 - 1} (\text{PRot}(d_j; -n_1 i) \odot \text{Rot}(ct_v; j \bmod n_1))$ for $i \in [0, n_2]$. ct_w denotes the ciphertext result of matrix-vector multiplication. $ct_w = \bigoplus_{i=0}^{n_2 - 1} \text{Rot}(d'_i; n_1 i)$. The number of rotations needed in BSGS diagonal method is determined by $n_1 = 8$ and $n_2 = 8$. It can be observed that $n_1 - 1 + n_2 - 1 = 14$ rotations are required. The number of plaintexts needed is still the same as the diagonal method. It needs $n = 64$ plaintexts.

The example of our proposed algorithm, parallel BSGS matrix-vector multiplication, can be seen in Figure 11. In Figure 11, we showcase the scenario where parallel BSGS matrix-vector multiplication is applied to the fully connected layer of Resnet (for CIFAR-10 Images) implementation of

[8]. In the input ciphertext ct_v , 8 identical data are located parallelly, and the length of the ciphertext is $n_t = 32768$. Unlike the conventional BSGS diagonal method, which only utilizes the leftmost data, we utilize all 8 data. To further clarify the concept, additional definitions are necessary to explain Figure 11. Here, p_i signifies that there are p_i identical data which positioned at equal intervals within a single ciphertext. Specifically, $p_i = 8$ in the case of ct_v . $(|v_0|v_1|v_2|v_3|v_4|v_5|v_6|v_7|)$ indicates that within one ciphertext, the value of v_i is located at positions $i n_t / 8$. D represents the arrangement of multiple d 's within one ciphertext. It can be precisely defined as follows:

$$D_i = \left(|d_i| \text{PRot}(d_{i+n_1}; -n_1) | \text{PRot}(d_{i+2n_1}; -2n_1) | \dots \right. \\ \left. \dots | \text{PRot}(d_{i+(n_2-1)n_1}; -(n_2-1)n_1) | \right) \quad (7)$$

for $0 \leq i < n_1$. In other words, D is obtained by arranging n_2 number of d 's within one ciphertext at equal intervals. D' represents a similar arrangement where multiple d 's exist within the ciphertext: $D' = \left(|d'_0|d'_1| \dots |d'_{n_2-1}| \right)$. As observed in Figure 11, D' can be obtained as the sum of the multiplications of D and the rotated ct_v . Ultimately, D' is generated through three rotations and additions to combine d' , thus reducing the number of rotations during the process of combining d' .

The key idea of the concept is to modify the plaintext such that, unlike the conventional structure where there was only one d within one ciphertext, multiple d 's exist, as illustrated by D . This idea increases the number of valid values within the plaintext and reduces the total number of plaintexts required. Additionally, by repetitively applying rotation and addition operations to the ciphertext after computation, d' can be effectively combined. In practice, the advantage is evident in the reduced number of rotations required, as observed in Figure 11 where the rotation is reduced to $2 \log_2 n_2 + n_1 - \log_2 p_i = 11$ where $n = 64 = n_1 n_2 = 8 \times 8$. Furthermore, the required number of plaintexts is reduced to $n_1 = 8$.

In Figure 11, using the structure where there are 8 identical data within the ciphertext remains unchanged during the process. However, increasing the number of identical data within the ciphertext as needed could be a method to reduce the number of rotation operations. For instance, let's consider increasing the number of identical data within the ciphertext to 16 by setting $n_1 = 4$ and $n_2 = 16$ and applying $ct_v = ct_v \oplus \text{Rot}(ct_v; 2048)$. In this scenario, when applying parallel BSGS matrix-vector multiplication, the required rotation count reduces to $2 \log_2 n_2 + n_1 - \log_2 p_i = 9$. The required number of plaintexts also decreases to $n_1 = 4$. Thus, depending on how n_1 and n_2 are set, performance can vary. Therefore, we provide an equation to determine n_2 and $n_1 = n/n_2$ for a given n , which minimizes the rotation when using parallel BSGS matrix-vector multiplication.

$$\underset{n_2}{\operatorname{argmin}} (2 \log_2 (n_2) + (n/n_2) - \log_2 p_i) \quad (8)$$

Algorithm 3 ParBSGSMatVecMul(D, ct_v, n_1, n_2, p_i)

Input: Plaintexts of parallelly located diagonals of $n \times n$ matrix D , Parallelly multiplexed tensor ciphertext that contains n valid value ct_v , number of identical data in ct_v p_i , n_1 and n_2 that satisfy $n = n_1 n_2$.

Output: Parallelly multiplexed tensor ciphertext that contains n valid value ct_w .

```

1:  $ct_w \leftarrow ct_{\text{zero}}$ 
2:  $ct_v \leftarrow ct_v \oplus \text{Rot}(ct_v; -n_1 n_2)$ 
3: for  $i_1 \leftarrow 1$  to  $\log_2(n_2/p_i)$  do
4:    $ct_v \leftarrow ct_v \oplus \text{Rot}(ct_v; -(n_t/p_i)/2^{i_1})$ 
5: end for
6: for  $i_1 \leftarrow 0$  to  $n_1 - 1$  do
7:    $ct_w \leftarrow ct_w \oplus D[i_1] \odot \text{Rot}(ct_v; i_1)$ 
8: end for
9: for  $i_1 \leftarrow 0$  to  $\log_2(n_2) - 1$  do
10:   $ct_w \leftarrow ct_w \oplus \text{Rot}(ct_w; 2^{i_1}(n_1 + n_t/n_2))$ 
11: end for
12: return  $ct_w$ 

```

At first glance, it may seem that increasing the value of n_2 , which determines the number of d in a single ciphertext, indefinitely could lead to more optimization. However, since the length of a single ciphertext is limited, the number of identical data that can exist within the ciphertext is also limited. In other words, in the Equation 8, n_2 is bounded as $p_i \leq n_2 \leq n_t/(2n)$. With the determined values of n_1 and n_2 obtained in this manner, we can employ our parallel BSGS matrix-vector multiplication (ParBSGSMatVecMul) algorithm, as detailed in Algorithm 3. ParBSGSMatVecMul's input, the parallel diagonal plaintext D can be further generalized as follows: $D_i = \sum_{j=0}^{n_2-1} \text{PRot}(d[jn_1+i]; -jn_1 - (n_t/n_2)j)$ where $i \in [0, n_1]$

Finally, in terms of guidance for practical usage, to effectively utilize ParBSGSMatVecMul, it is necessary for there to be multiple identical data within a single ciphertext. Therefore, to use ParBSGSMatVecMul, 1) the number of slots for ciphertexts and plaintexts should be large relative to the size of the matrix and vector, and 2) multiple instances of the same data should be allowed within a ciphertext. Additionally, if there are already multiple identical data within a ciphertext (denoted as p_i), the effectiveness of ParBSGSMatVecMul can be enhanced, although the value of p_i is not critical to performance.

6. Rotation Key Reduction

Rotation keys required by a given FHE algorithm are typically generated and transmitted from the client to the server prior to execution, and are then stored for use during computation. Therefore, they do not directly affect the latency of the algorithm itself. However, if the number of such required rotation keys is excessively large, it may result in substantial transmission latency, which can become a practical bottleneck in real-world deployments.

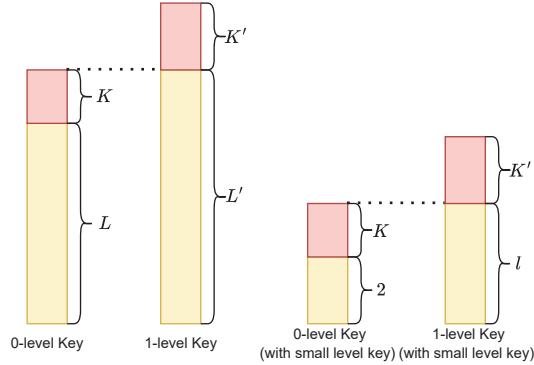


Figure 12: A simple and abstract illustration of applying the hierarchical rotation key system and small-level key system to rotation key.

While proposing RO-MPConv, we significantly reduced the number of rotations and effectively decreased the number of required rotation keys by about one-third. In this section, we aim to further minimize the size of the rotation keys that need to be transmitted from the client to the server using additional methods. We employ the hierarchical rotation key system introduced by [11]. The hierarchical rotation key system organizes computation keys into hierarchical levels, reducing the size of the keys that need to be transmitted from the client to the server. If we define the keys required to perform rotations on the server as 0 key-level rotation keys, the client only sends 1 key-level rotation keys. Then, the server generates 0 key-level rotation keys from the 1 key-level rotation keys. This system, known as the two-level hierarchical rotation key system, allowed us to reduce the size of rotation keys.

Furthermore, additional optimization is possible by leveraging the structure of rotation keys. In CKKS, rotation keys have a combined structure supporting rotation operations from level 0 up to the maximum mult level in CKKS parameters to ensure support for rotation operations at all mult levels. However, if we know precisely at which mult level rotation operations occur, we can transmit only the keys relevant to that mult level, reducing the transmission volume. When applied to RO-MPConv technology, by default, 2 levels are consumed during RO-MPConv operations. Therefore, transmitting only the rotation keys corresponding to these levels can reduce the transmission volume of rotation keys. We define the technique that reduces the size of the rotation key by limiting the rotation key to support rotation operations at a limited level as *small-level key system*. We ultimately extended this to a two-level hierarchical rotation key system.

A simple illustration of the key structure used in this method is shown in Figure 12. 0-level key in Figure 12 represents conventional rotation key. L denotes the maximum level that the ciphertext can have, and K corresponds to the special modulus. In our implementation environment, L and K are 24 and 5, respectively. By applying the two-level hierarchical rotation key system, we can derive the 1-level

Table 2

Comparison of rotation key size and transmission latency, assuming 100Mbps for network connection. The evaluation is conducted on all the convolution operations utilizing 2 depth used in ResNet that take CIFAR-10 images as input.

Method	MP-Conv [8]	RO-MPConv	RO-MPConv +Hierarchical rotation key system	RO-MPConv +Hierarchical rotation key system +Small-level key system
Rotation key size(MB)	29100	8250	2100	1100
Transmission latency(min)	38.8	11.0	2.8	1.5

key, whose maximum level becomes $L' = L + K = 29$. The key labeled as “1-level Key” in the figure corresponds to this. Furthermore, by applying our small-level key system, the supported level of the 0-level key can be reduced to 2. As a result, the maximum level of the 1-level key is also significantly reduced to $l = K+2 = 7$, allowing us to achieve a significant reduction in the size of each rotation key.

Table 2 represents the experiment results of applying a two-level hierarchical rotation key system and small-level key system to reduce the size of rotation keys required for RO-MPConv. As mentioned in Subsection 4.3, utilizing additional depth is unrelated to the number of rotation keys. Therefore, experiments were conducted using the convolutions utilizing 2 depth. By applying the hierarchical rotation key system and small-level key system to the rotation keys transmitted in RO-MPConv, we were able to reduce the size of the final transmitted rotation keys by approximately 29× compared to [8]. Transmission latency showed similar results, decreasing from 38.8 minutes to 1.5 minutes, a result that suggests a meaningful improvement in the method’s practical applicability.

7. Experimental Results

In this section, we will compare the execution time and size of rotation keys of the various operations proposed by us with those of the conventional method. The numerical analyses are conducted on the representative Lattigo v5 (lattice-based homomorphic encryption library, [40]) on AMD Ryzen 9 7900X at 3.2 GHz (24 cores) with 124 GB RAM, running the Ubuntu 20.04 operating system. We set Lattigo main CKKS parameter as follows: $N = 2^{16}$, $\log Q = [51, 46]$, $\log P = [60, 60, 60, 60, 60]$. P denotes the special modulus, and length of special modulus correspondence to the length of $\log P$. Since our algorithm relies solely on standard CKKS operations, it does not exhibit any particular parameter settings under which it performs exceptionally better or worse. To ensure the reproducibility of our results, the complete source code and

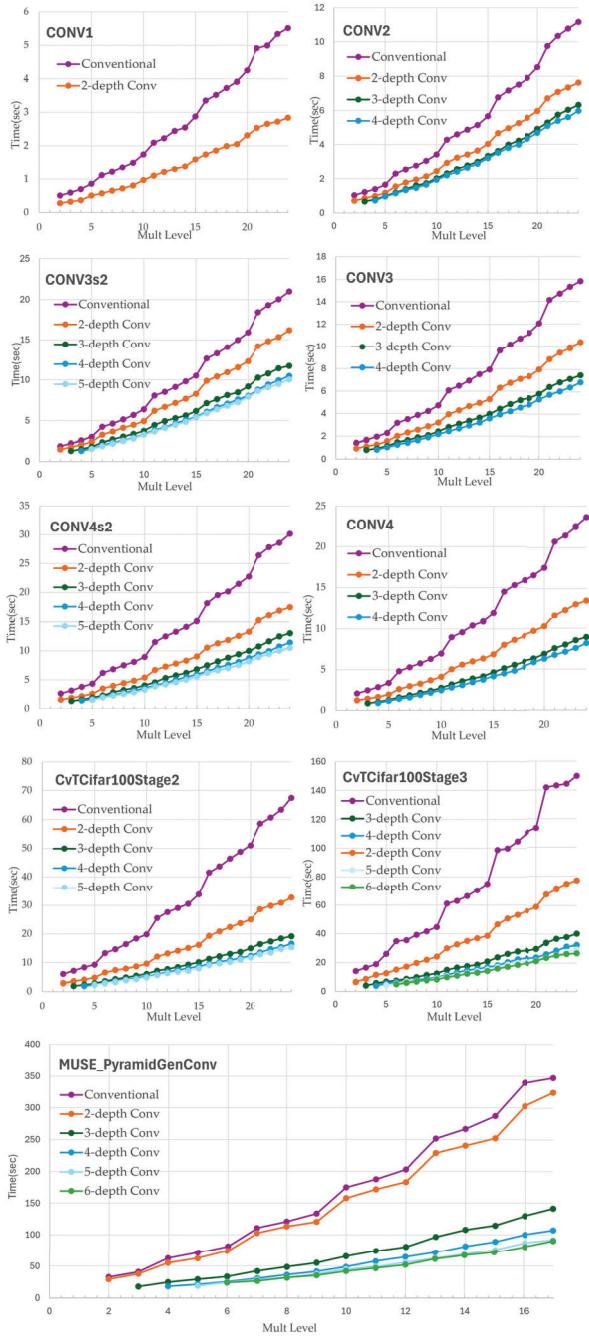


Figure 13: Latency comparison of proposed (RO-MPConv) and conventional (multiplexed parallel convolution in [8]) convolution.

scripts for all convolution and matrix-vector multiplication experiments presented in this paper have been made available at our GitHub repository

To ensure a comprehensive evaluation of RO-MPConv, we benchmarked our algorithms on nine distinct convolution configurations. Six of these (CONV1, CONV2, CONV3s2, CONV3, CONV4s2, CONV4) are based on the ResNet architecture for CIFAR-10, following the precedent set in the

original MPConv study. To demonstrate broader applicability, we also implemented three convolution operations from modern, state-of-the-art AI models: two from the Convolutional Vision Transformer (CvT) [4] for CIFAR-100 inference (CvTConv1, CvTConv2) and one from the MUSE [7] model's multi-scale feature pyramid generator (MuseConv). The detailed parameters for these convolutions, such as their tensor shapes and stride values, are available in Appendix B. In the following discussion and figures, an i -depth Conv refers to a RO-MPConv variant specifically designed to consume i multiplicative levels during its execution.

Figure 13 presents the latency comparison of our RO-MPConv against the conventional MPConv. It should be noted that for MuseConv, due to memory constraints caused by the massive rotation key requirements of the conventional MPConv, we conducted experiments for MuseConv using a lighter CKKS parameter setting with a maximum level of 18. The results of Figure 13 clearly show that RO-MPConv significantly reduces execution time across all tested configurations, with performance gains becoming more pronounced at higher multiplicative levels. For the ResNet-based convolutions, our 2-depth RO-MPConv achieves latency reductions of 20% to 45%, with further improvements up to 66%, 70% when using additional depth (3-depth Conv, 4-depth Conv). The performance gains were even more substantial for modern architectures. For the CvT operations, the 2-depth variant reduced latency by approximately 56%, while the 3-depth and 4-depth variants achieved reductions of up to 79% and 81%, respectively. In the case of MuseConv, the 2-depth implementation offered a modest 10% improvement, but the 3-depth and 4-depth versions provided 62%, 72% latency reduction.

The source of these improvements is detailed in Table 3 which quantifies the reduction in the number of rotation operations and rotation keys. The most significant reduction occurs in the ZeroOutCombine process, where our novel CrossCombine subprocess replaces a number of rotations proportional to the output channels (c_o) with a much smaller number related to the parallel data count ($p_c - 1$). This reduction is particularly impactful for convolutions with a large number of channels, such as MuseConv. Although its 2-depth latency improvement was moderate, Table 3 reveals advantage that RO-MPConv reduced the number of required rotation keys for MuseConv by approximately 95%.

To demonstrate the practical impact of these algorithmic improvements on end-to-end FHE-based PPML inference, we integrated RO-MPConv into state-of-the-art models that originally utilized MPConv. We selected representative models from three distinct categories to showcase the broad applicability and effectiveness of our approach. These categories include: 1) **structurally modified CNNs** such as the Patch CNN (PCNN) in CryptoFace [10], which alters the standard sequential architecture to manage depth consumption in face recognition; 2) **mixed-degree CNNs** from AutoFHE [9], which employs varying degrees of polynomial approximations for ReLU across different layers; and 3) **low-degree approximations**, represented by a baseline

Table 3

Comparison of the number of rotations and rotation keys used in the RotationSum and ZeroOutCombine processes of convolution. The SISOConv process, which consistently requires 8 rotations and 8 rotation keys and remains unchanged between the conventional and proposed methods, is omitted to reduce redundancy. Total = SISOConv (8) + RotationSum + ZeroOutCombine.

Convolution	Process	MPConv [8]	No. of rotations						No. of rotation keys	
			2-depth Conv	3-depth Conv	4-depth Conv	5-depth Conv	6-depth Conv	MPConv [8]	2-depth Conv	
CONV1	RotationSum	4	4	-	-	-	-	2	3	
	ZeroOutCombine	17	2	-	-	-	-	16	2	
	Total	29	14	-	-	-	-	26	13	
CONV2	RotationSum	32	24	16	14	-	-	4	6	
	ZeroOutCombine	17	4	4	4	-	-	16	4	
	Total	57	36	28	26	-	-	28	18	
CONV3s2	RotationSum	64	64	40	32	30	-	4	8	
	ZeroOutCombine	34	5	5	5	5	-	33	5	
	Total	106	77	53	45	43	-	45	21	
CONV3	RotationSum	40	32	22	18	-	-	5	7	
	ZeroOutCombine	34	9	5	5	-	-	33	9	
	Total	82	49	35	31	-	-	46	24	
CONV4s2	RotationSum	80	64	44	36	32	-	5	8	
	ZeroOutCombine	67	10	6	6	6	-	66	10	
	Total	155	82	58	50	46	-	79	26	
CONV4	RotationSum	48	48	24	20	-	-	6	9	
	ZeroOutCombine	67	10	10	10	-	-	66	10	
	Total	123	66	42	38	-	-	80	27	
CvTConv1	RotationSum	144	120	72	57	51	-	6	10	
	ZeroOutCombine	195	18	10	10	10	-	194	18	
	Total	347	146	90	75	69	-	208	36	
CvTConv2	RotationSum	384	288	156	126	108	102	8	12	
	ZeroOutCombine	388	35	19	11	11	11	387	35	
	Total	780	331	183	145	127	121	403	55	
MuseConv	RotationSum	4608	4608	1920	1408	1184	1088	9	18	
	ZeroOutCombine	512	0	0	0	0	0	511	0	
	Total	5128	4616	1928	1416	1192	1096	528	26	

configuration utilized in both the AutoFHE and CryptoFace studies, where the MP-CNN activation function is replaced with a low-degree polynomial from AESPA [41] to prioritize speed over accuracy.

MPConv's model-agnostic nature makes it a straightforward choice for these diverse applications. However, these advanced models often incur higher convolution-related overhead compared to the original MP-CNN, a challenge explicitly reported in studies like AutoFHE and CryptoFace. This trend is corroborated by our findings in Table 4, where the "Convolution latency fraction" for these models is notably higher than that of the baseline MP-CNN.

The primary reason for this increased overhead is illustrated in Figure 14. The original MP-CNN architecture was specifically designed to execute all MPConv operations at the lowest possible multiplicative levels, where they are fastest (Figure 14a). In contrast, models like AutoFHE introduce structural variations and mixed-degree activation functions, which force convolution operations to be executed

at higher, and thus slower, levels (Figure 14b). This often leaves small, unused gaps in the level consumption between layers.

Our solution directly targets this inefficiency. As shown in Figure 14c, we not only replace MPConv with the fundamentally faster RO-MPConv but also strategically utilize RO-MPConv variants that consume additional depth. This allows us to perfectly fill the previously wasted multiplicative levels, creating a more compact and efficient computational graph. Priority for using these depth-consuming variants is given to convolutions operating at higher levels, where the potential latency reduction is greatest.

By applying this substitution strategy, we achieved significant performance gains. As detailed in Table 4, the latency of the convolution modules themselves was reduced by up to **57.2%**. More importantly, this translated into a substantial reduction in the total end-to-end model latency by up to **26.04%**. While RO-MPConv also reduces convolution latency in the baseline MP-CNN, the overall improvement is

Table 4

Performance improvements from replacing MPConv with RO-MPConv in state-of-the-art HE-based models. Consistent with the experimental settings in AutoFHE and CryptoFace, these evaluations were conducted using a maximum multiplicative level of 30, with bootstrapping consuming 14 levels. The table shows the fraction of total latency attributed to convolution operations, the latency improvement of the convolution module itself, and the resulting improvement in total end-to-end inference latency. The results demonstrate that RO-MPConv provides substantial system-level performance gains, particularly for models with a high convolution latency fraction.

	Model	Convolution latency fraction	MPConv (sec)	RO-MPConv (sec)	Convolution latency improvement	Total latency improvement
MP-CNN	ResNet20	15.2%	31.89	20.17	36.8%	5.59%
	ResNet32	14.7%	49.95	31.69	36.6%	5.37%
	ResNet44	14.3%	68.02	43.21	36.5%	5.22%
AESPA	ResNet20	52.0%	120.64	75.75	37.2%	19.36%
	ResNet32	52.0%	190.81	119.67	37.3%	19.37%
	ResNet44	51.0%	260.98	163.59	37.3%	19.02%
AutoFHE	ResNet20 (boot5)	47.7%	96.93	51.37	47.0%	22.43%
	ResNet20 (boot11)	30.6%	102.16	44.7	56.2%	17.20%
	ResNet32 (boot8)	53.0%	155.27	95.24	38.7%	20.47%
	ResNet32 (boot19)	25.1%	125.84	55.46	55.9%	14.05%
	ResNet44 (boot8)	58.1%	185.6	102.48	44.8%	26.04%
	ResNet44 (boot22)	28.7%	156.92	67.21	57.2%	16.39%
CryptoFace	CryptoFace4	62.9%	305.39	199.18	34.8%	21.89%

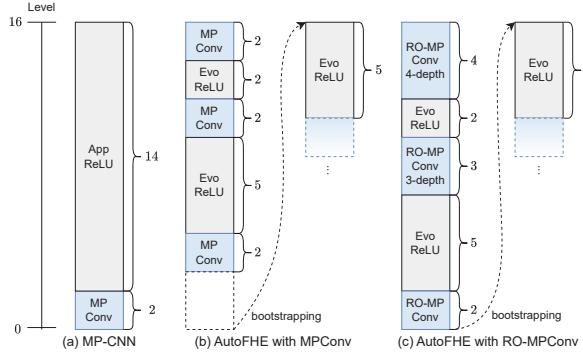


Figure 14: Comparison of operational layouts across multiplicative levels in different model architectures. (a) The baseline MP-CNN, which utilizes a deep polynomial activation (AppReLU), requiring MPConv operations to be scheduled at the lowest (fastest) levels immediately before bootstrapping. (b) AutoFHE with standard MPConv, where the use of mixed-degree activations (EvoReLU) results in convolutions being scheduled at higher levels and creates intermittent unused level gaps. (c) Our proposed application, where MPConv is replaced by RO-MPConv. By employing RO-MPConv variants with additional depth (e.g., 3-depth, 4-depth), we can efficiently schedule operations to fill these previously wasted levels, creating an optimized computational graph and reducing overall latency.

modest because the convolution fraction was already small. However, in advanced models where architectural trade-offs pushed convolution to higher levels and increased its latency contribution, RO-MPConv proved to be highly effective at reducing the entire model’s inference time.

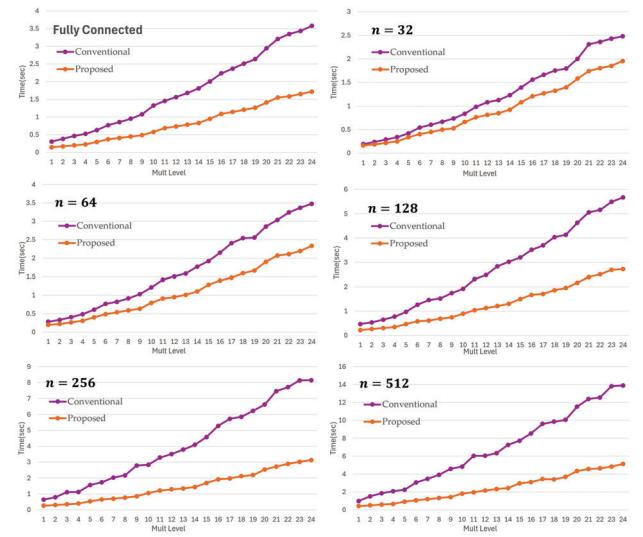


Figure 15: Latency comparison of parallel BSGS matrix-vector multiplication and BSGS diagonal method.

As an additional contribution, this paper introduces a parallel Baby-Step Giant-Step (BSGS) matrix-vector multiplication algorithm designed to accelerate other critical linear operations, such as fully connected layers. The effectiveness of this method is demonstrated in Figure 15, which compares its latency against the conventional BSGS diagonal method [12] under various scenarios.

First, when applied to the fully connected layer configuration from [8] (where $n = 64$, $p_i = 8$, and $n_t = 32768$), our proposed method achieves a nearly 2x speedup across most

Table 5

Comparison of the number of rotations and plaintexts required for matrix-vector multiplication between the BSGS diagonal method and the proposed ParBSGSMatVecMul algorithm. The scenario is same as Figure 15. Assume matrix size as $n \times n$, and vector size as $n \times 1$.

Scenario	Criteria	BSGS diagonal method	Proposed
Fully connected layer in [8]	# Rotation	15	9
	# Plaintext	64	2
$n = 32$	# Rotation	11	10
	# Plaintext	32	2
$n = 64$	# Rotation	15	12
	# Plaintext	64	2
$n = 128$	# Rotation	23	14
	# Plaintext	128	2
$n = 256$	# Rotation	31	16
	# Plaintext	256	4
$n = 512$	# Rotation	47	26
	# Plaintext	512	16

multiplicative levels. Second, to assess its general applicability, we evaluated its performance on $n \times n$ matrix-vector multiplications for various n , setting $p_i = 1$ to represent cases without pre-existing parallel data. Even in this more challenging setting, our algorithm achieved a substantial latency reduction of up to 69%.

This performance gain stems from a fundamental reduction in the required number of rotation operations. For a matrix of size $n = n_1 n_2$, the conventional BSGS diagonal method requires approximately $n_1 + n_2 - 1$ rotations. In contrast, our parallel approach requires only $2 \log_2(n_2) + n_1 - \log_2 p_i$ rotations. Due to the logarithmic term in our formula, the efficiency advantage of our method becomes more pronounced as the matrix size n increases. This theoretical property is clearly reflected in the performance trends shown in Figure 15. Table 5 provides a detailed breakdown of the number of rotations and plaintexts required for each scenario, further quantifying the improvements.

8. Implementation Considerations

Our experimental results have demonstrated that RO-MPConv is a versatile optimization, successfully enhancing not only ResNet but also a range of modern SOTA architectures like CVT and MUSE. However, similar to all algorithms built upon FHE, its implementation operates within the fundamental constraint of a finite ciphertext length. This length is a critical system parameter, as arbitrarily increasing it to accommodate larger inputs would introduce prohibitive computational latency.

Consequently, for any FHE-based algorithm, a standard approach is required when input data dimensions exceed the capacity of a single ciphertext: partitioning the data across

multiple ciphertexts. To provide clear guidance for implementation, we explicitly define the operational boundaries of our algorithm, preempting potential integration challenges.

The following equations formalize the precise conditions under which the proposed RO-MPConv can be applied directly without requiring such data partitioning.

$$2^{\lceil \log_2 h_i \rceil + \lceil \log_2 w_i \rceil} c_i \leq n_t \quad (9)$$

$$2^x c_o \leq 2^{\lceil \log_2 c_i \rceil} p_i \quad (10)$$

The Equation 9 expresses the constraint that the input size of the convolution must not exceed the length of the ciphertext. The Equation 10 addresses whether the ciphertext length is sufficient to generate the convolution output. Specifically, p_i represents the number of identical data elements in the ciphertext, and each data element can hold $2^{\lceil \log_2 c_i \rceil}$ channels of information. If the required number of output channels exceeds this capacity, the ciphertext length would be insufficient. Thus, the relationship $c_o \leq 2^{\lceil \log_2 c_i \rceil} p_i$ is defined. Moreover, when considering the correlation between the RotationSum process and the CrossCombine process in Section 4.2, the number of data combining during CrossCombine increases by a factor of 2^x . This leads to the final equation $2^x c_o \leq 2^{\lceil \log_2 c_i \rceil} p_i$. In addition, by rearranging Equation 10 with respect to x , the upper bound for x can be directly obtained as $x \leq \log_2(2^{\lceil \log_2 c_i \rceil} p_i / c_o)$.

If a target convolution does not satisfy these conditions, the input or output tensors must be partitioned before applying our algorithm to each segment. This delineation provides a practical guide for integrating RO-MPConv into future HE-based systems.

9. Conclusion and Future Works

In this paper, we addressed the critical performance bottlenecks that hinder the practical deployment of homomorphic encryption-based deep learning: high latency and large key sizes, particularly in convolution operations. We introduced Rotation-Optimized Multiplexed Parallel Convolution (RO-MPConv), a novel algorithm that fundamentally redesigns the state-of-the-art MPConv to minimize the number of costly rotation operations. We further enhanced its flexibility by introducing variants that consume additional multiplicative depth, offering a strategic trade-off for even greater latency reductions.

Our experimental results validate the effectiveness of these contributions. RO-MPConv not only significantly reduces the latency of standalone convolution operations but, more importantly, improves the end-to-end inference performance of existing state-of-the-art models by up to 26% through direct substitution. This demonstrates that optimizing core HE primitives translates into substantial system-level gains. Furthermore, by combining our proposed small-level key system with a hierarchical structure, we reduced

the rotation key transmission size by a factor of 29, directly tackling the communication overhead in client-server environments. As a complementary contribution, we also presented a parallel BSGS matrix-vector multiplication method, which accelerates other common linear operations by leveraging the data packing characteristics inherent in many HE schemes.

Collectively, our work provides a suite of optimizations that significantly lowers the computational and communication barriers for privacy-preserving machine learning. By making fundamental operations faster and lighter, we enhance the practical feasibility of deploying complex, deep neural networks over encrypted data. Building on this foundation, our future work will focus on utilizing these optimized building blocks to construct fully homomorphic implementations of modern, convolution-integrated architectures, such as advanced CNN variants, convolution-integrated Vision Transformers, and State-Space Models, further expanding the horizons of secure and private AI applications.

A. RO-MPConv blueprints

This section provides the **blueprints** used for implementing 2-depth Conv, 3-depth Conv, 4-depth Conv, and 5-depth Conv for the six RO-MPConvs utilized in ResNet architectures. While this paper includes a substantial portion of the blueprints for rotation-optimized convolutions, it does not cover all cases due to the large number of scenarios explored. However, all **blueprints** for the RO-MPConvs we have conducted can be accessed via our GitHub link.

All **blueprints** are designed with consideration for the relationship between the RotationSum and CrossCombine process in RO-MPConv, and all **blueprints** are designed with the goal of minimizing time. These blueprints may not be the only blueprints in each convolution scenario. There could be other blueprints with different values, but they are unlikely to have better performance than our blueprint. In other words, there are many correct blueprints, but not all of them are optimal. Various other correct blueprints can also be tested through our GitHub link.

RO-MPConv has some heuristics nature, so here's a simple guide for creating **blueprints**. Firstly, rather than focusing on the correlation between processes RotationSum and CrossCombine, create the **blueprint** backward from the final output of the convolution. In other words, decide how to divide the final output of the convolution for CrossCombine, then determine the order for performing RotationSum to move each channel to its position, and finally, decide on the order of kernel weights. This approach makes it easier to create the initial **blueprint**. Subsequently, you can further optimize the **blueprint** using the correlation between processes RotationSum and CrossCombine. Additionally, if utilizing additional depth is allowed, you can combine ciphertext that requires the same shift of rotation during RotationSum for further optimization on execution time.

It should be noted that in all expressions related to **blueprint** algorithms, the x -th element of a 1-dimensional **blueprint** A is denoted as $A[x]$, and for some 2-dimensional **blueprint** B, the y -th element of $B[x]$ is denoted as $B[x][y]$. The index of the first value of **blueprint** is 0, for instance, the first value of 1-dimensional **blueprint** A is $A[0]$. KernelBP is set to be consistent regardless of how many depths the convolution uses, for the sake of simplicity.

B. Parameters

Various parameters such as $h_i, h_o, w_i, w_o, c_i, c_o, f_h, f_w, s, k_i, k_o, t_i, t_o, p_i, p_o, p_c$ and q are defined differently depending on the shape of each convolution. Table 11 shows the values of these parameters used in each shape of convolution. These parameters are also defined to represent the MPConv described in [8], to avoid confusion, a similar notation as in [8] was used.

Parameters w_i, h_i , and c_i respectively represent the width, height, and channel numbers of the three-dimensional input data of convolution. Similarly, w_o, h_o , and c_o respectively represent the width, height, and channel numbers of the three-dimensional output data of convolution. f_h, f_w are height and width numbers of kernel data, and s represents stride value of convolution. k_i, k_o are gaps of input and output ciphertext of convolution, respectively. Other variable can be defined as $t_i = \lceil \frac{c_i}{k_i^2} \rceil$, $t_o = \lceil \frac{c_o}{k_o^2} \rceil$, $p_i = 2^{\lfloor \log_2 \left(\frac{n_t}{k_i^2 h_i w_i t_i} \right) \rfloor}$, $p_o = 2^{\lfloor \log_2 \left(\frac{n_t}{k_o^2 h_o w_o t_o} \right) \rfloor}$, and $q = \lceil \frac{c_o}{p_i} \rceil$. Specifically, p_i and p_o represent the number of identical parallelly existing data within input ciphertext, and output ciphertext, respectively. p_c is a new parameter proposed to denote the essential number of parallel data to be merged in CrossCombine. $p_c = \max(s^2, p_i / \lceil p_i / p_o \rceil)$.

Table 6
5-depth Conv **blueprint**.

Convolution	RotationSumBP	CrossCombineBP
CONV3s2	[[4], [1, 2, 1024], [1, 2, 2048], [1, 2, 4096], [2, 2, 8192]]	[4, 0, 8191, 16352, 24543]
CONV4s2	[[5], [1, 2, 1], [1, 2, 32], [1, 2, 1024], [0, 4096], [2, 2, 2048]]	[4, 0, 8190, 16320, 24510]

Table 7
4-depth Conv **blueprint**.

Convolution	RotationSumBP	CrossCombineBP
CONV2	[[3], [1, 2, 1024], [1, 2, 2048], [2, 2, 4096]]	[4, 0, 8192, 8192, 16384]
CONV3s2	[[3], [1, 2, 1024], [1, 2, 2048], [2, 4, 4096, 8192]]	[4, 0, 8191, 16352, 24543]
CONV3	[[5], [1, 2, 1], [1, 2, 32], [0, 2048], [0, 4096], [2, 2, 1024]]	[4, 0, 6144, 12288, 18432]
CONV4s2	[[4], [1, 2, 1], [1, 2, 32], [0, 4096], [2, 4, 1024, 2048]]	[4, 0, 8190, 16320, 24510]
CONV4	[[6], [1, 2, 1], [1, 2, 2], [0, 64], [0, 1024], [0, 2048], [2, 2, 32]]	[8, 0, 4032, 7168, 11200, 14336, 18368, 21504, 25536]

Table 8
3-depth Conv **blueprint**.

Convolution	RotationSumBP	CrossCombineBP
CONV2	[[2], [1, 2, 1024], [2, 4, 2048, 4096]]	[4, 0, 8192, 8192, 16384]
CONV3s2	[[2], [1, 4, 1024, 2048], [2, 4, 4096, 8192]]	[4, 0, 8191, 16352, 24543]
CONV3	[[4], [1, 4, 1, 32], [0, 2048], [0, 4096], [2, 2, 1024]]	[4, 0, 6144, 12288, 18432]
CONV4s2	[[3], [1, 4, 1, 32], [0, 4096], [2, 4, 1024, 2048]]	[4, 0, 8190, 16320, 24510]
CONV4	[[5], [1, 4, 1, 2], [0, 64], [0, 1024], [0, 2048], [2, 2, 32]]	[8, 0, 4032, 7168, 11200, 14336, 18368, 21504, 25536]

Table 9
2-depth Conv **blueprint**.

Convolution	RotationSumBP	CrossCombineBP
CONV1	[[2], [0, 2048], [2, 2, 1024]]	[2, 0, 14336]
CONV2	[[1], [2, 8, 1024, 2048, 4096]]	[4, 0, 8192, 8192, 16384]
CONV3s2	[[1], [2, 16, 1024, 2048, 4096, 8192]]	[4, 0, 8191, 16352, 24543]
CONV3	[[2], [0, 2048], [2, 8, 1, 32, 1024]]	[8, 0, 4096, 6144, 10240, 12288, 16384, 18432, 22528]
CONV4s2	[[1], [2, 16, 1, 32, 1024, 2048]]	[8, 0, 4096, 8190, 12286, 16320, 20416, 24510, 28606]
CONV4	[[4], [0, 64], [0, 1024], [0, 2048], [2, 8, 1, 2, 32]]	[8, 0, 4032, 7168, 11200, 14336, 18368, 21504, 25536]

Table 10
KernelBP **blueprint**.

Convolution	KernelBP
CONV1	[[0, 4, 8, 12, 2, 6, 10, 14], [1, 5, 9, 13, 3, 7, 11, 15]]
CONV2	[[0, 8], [1, 9], [2, 10], [3, 11], [4, 12], [5, 13], [6, 14], [7, 15]]
CONV3s2	[[0, 2], [4, 6], [8, 10], [12, 14], [16, 18], [20, 22], [24, 26], [28, 30], [1, 3], [5, 7], [9, 11], [13, 15], [17, 19], [21, 23], [25, 27], [29, 31]]
CONV3	[[0, 8, 16, 24], [1, 9, 17, 25], [2, 10, 18, 26], [3, 11, 19, 27], [4, 12, 20, 28], [5, 13, 21, 29], [6, 14, 22, 30], [7, 15, 23, 31]]
CONV4s2	[[0, 2, 8, 10], [1, 3, 9, 11], [4, 6, 12, 14], [5, 7, 13, 15], [16, 18, 24, 26], [17, 19, 25, 27], [20, 22, 28, 30], [21, 23, 29, 31], [32, 34, 40, 42], [33, 35, 41, 43], [36, 38, 44, 46], [37, 39, 45, 47], [48, 50, 56, 58], [49, 51, 57, 59], [52, 54, 60, 62], [53, 55, 61, 63]]
CONV4	[[0, 8, 16, 24, 32, 40, 48, 56], [1, 9, 17, 25, 33, 41, 49, 57], [2, 10, 18, 26, 34, 42, 50, 58], [3, 11, 19, 27, 35, 43, 51, 59], [4, 12, 20, 28, 36, 44, 52, 60], [5, 13, 21, 29, 37, 45, 53, 61], [6, 14, 22, 30, 38, 46, 54, 62], [7, 15, 23, 31, 39, 47, 55, 63]]

Table 11

Parameters that are used in MPConv and RO-MPConv.

Component	CONV1	CONV2	CONV3s2	CONV3	CONV4s2	CONV4	CvTConv1	CvTConv2	MuseConv
f_h	3	3	3	3	3	3	3	3	3
f_w	3	3	3	3	3	3	3	3	3
s	1	1	2	1	2	1	2	2	1
h_i	32	32	32	16	16	8	8	4	8
h_o	32	32	16	16	8	8	4	2	8
w_i	32	32	32	16	16	8	8	4	8
w_o	32	32	16	16	8	8	4	2	8
c_i	3	16	16	32	32	64	64	192	512
c_o	16	16	32	32	64	64	192	384	512
k_i	1	1	1	2	2	4	1	2	1
k_o	1	1	2	2	4	4	2	4	1
t_i	3	16	16	8	8	4	64	48	512
t_o	16	16	8	8	4	4	48	24	512
p_i	8	2	2	4	4	8	8	8	1
p_o	2	2	4	4	8	8	8	16	1
p_c	2	2	4	4	4	8	8	8	1
q	2	8	16	8	16	8	24	48	512

CRediT authorship contribution statement

Byeong-Seo Min: Writing – review & editing, Writing – original draft, Visualization, Software, Methodology, Investigation, Formal analysis. **Joon-Woo Lee:** Writing – review & editing, Validation, Supervision, Resources, Funding acquisition, Conceptualization.

Data Availability Statement

The source code developed for this study is openly available in a GitHub repository at <https://github.com/byeongseomin51/RO-MPConv.git>.

References

- [1] R. Xu, N. Baracaldo, J. Joshi, Privacy-preserving machine learning: Methods, challenges and directions, arXiv preprint arXiv:2108.04417 (2021).
- [2] J. H. Cheon, K. Han, A. Kim, M. Kim, Y. Song, A full rns variant of approximate homomorphic encryption, in: Selected Areas in Cryptography–SAC 2018: 25th International Conference, Calgary, AB, Canada, August 15–17, 2018, Revised Selected Papers 25, Springer, 2019, pp. 347–368.
- [3] J. H. Cheon, A. Kim, M. Kim, Y. Song, Homomorphic encryption for arithmetic of approximate numbers, in: Advances in Cryptology–ASIACRYPT 2017: 23rd International Conference on the Theory and Applications of Cryptology and Information Security, Hong Kong, China, December 3–7, 2017, Proceedings, Part I 23, Springer, 2017, pp. 409–437.
- [4] H. Wu, B. Xiao, N. Codella, M. Liu, X. Dai, L. Yuan, L. Zhang, Cvt: Introducing convolutions to vision transformers, in: Proceedings of the IEEE/CVF international conference on computer vision, 2021, pp. 22–31.
- [5] A. Gulati, J. Qin, C.-C. Chiu, N. Parmar, Y. Zhang, J. Yu, W. Han, S. Wang, Z. Zhang, Y. Wu, et al., Conformer: Convolution-augmented transformer for speech recognition, arXiv preprint arXiv:2005.08100 (2020).
- [6] H. Yan, Z. Li, W. Li, C. Wang, M. Wu, C. Zhang, Contnet: Why not use convolution and transformer at the same time?, arXiv preprint arXiv:2104.13497 (2021).
- [7] H. Tang, M. Cao, J. Huang, R. Liu, P. Jin, G. Li, X. Liang, Muse: Mamba is efficient multi-scale learner for text-video retrieval, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 39, 2025, pp. 7238–7246.
- [8] E. Lee, J.-W. Lee, J. Lee, Y.-S. Kim, Y. Kim, J.-S. No, W. Choi, Low-complexity deep convolutional neural networks on fully homomorphic encryption using multiplexed parallel convolutions, in: International Conference on Machine Learning, PMLR, 2022, pp. 12403–12422.
- [9] W. Ao, V. N. Boddeti, {AutoFHE}: Automated adaption of {CNNs} for efficient evaluation over {FHE}, in: 33rd USENIX Security Symposium (USENIX Security 24), 2024, pp. 2173–2190.
- [10] W. Ao, V. N. Boddeti, Cryptoface: End-to-end encrypted face recognition, in: Proceedings of the Computer Vision and Pattern Recognition Conference, 2025, pp. 19197–19206.
- [11] J.-W. Lee, E. Lee, Y.-S. Kim, J.-S. No, Rotation key reduction for client-server systems of deep neural network on fully homomorphic encryption, in: International Conference on the Theory and Application of Cryptology and Information Security, Springer, 2023, pp. 36–68.
- [12] S. Halevi, V. Shoup, Algorithms in helib, in: Advances in Cryptology–CRYPTO 2014: 34th Annual Cryptology Conference, Santa Barbara, CA, USA, August 17–21, 2014, Proceedings, Part I 34, Springer, 2014, pp. 554–571.
- [13] J. Wan, D. Li, J. Fang, Z. L. Jiang, Lpfhe: Low-complexity polynomial cnns for secure inference over fhe, in: European Symposium on Research in Computer Security, Springer, 2024, pp. 403–423.
- [14] N. Samardzic, D. Sanchez, Bitpacker: Enabling high arithmetic efficiency in fully homomorphic encryption accelerators, in: Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2, 2024, pp. 137–150.
- [15] R. Agrawal, A. Chandrakasan, A. Joshi, Heap: A fully homomorphic encryption accelerator with parallelized bootstrapping, in: 2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA), IEEE, 2024, pp. 756–769.
- [16] X. Deng, S. Fan, Z. Hu, Z. Tian, Z. Yang, J. Yu, D. Cao, D. Meng, R. Hou, M. Li, et al., Trinity: A general purpose fhe accelerator, in: 2024 57th IEEE/ACM International Symposium on Microarchitecture (MICRO), IEEE, 2024, pp. 338–351.
- [17] S. Cheon, Y. Lee, D. Kim, J. M. Lee, S. Jung, T. Kim, D. Lee, H. Kim, {DaCapo}: Automatic bootstrapping management for efficient fully homomorphic encryption, in: 33rd USENIX Security Symposium

- (USENIX Security 24), 2024, pp. 6993–7010.
- [18] D. Kim, J. Park, J. Kim, S. Kim, J. H. Ahn, Hyphen: A hybrid packing method and its optimizations for homomorphic encryption-based neural networks, *IEEE Access* (2023).
 - [19] H. Peng, R. Ran, Y. Luo, J. Zhao, S. Huang, K. Thorat, T. Geng, C. Wang, X. Xu, W. Wen, et al., Lingcn: Structural linearized graph convolutional network for homomorphically encrypted inference, *Advances in Neural Information Processing Systems* 36 (2024).
 - [20] R. Ran, X. Luo, W. Wang, T. Liu, G. Quan, X. Xu, C. Ding, W. Wen, Spencnn: orchestrating encoding and sparsity for fast homomorphically encrypted neural network inference, in: *International Conference on Machine Learning*, PMLR, 2023, pp. 28718–28728.
 - [21] B. Yalavarthi, C. Jutla, N. Ratha, Efficient convolution operator in fhe using summed area table, in: *International Conference on Pattern Recognition*, Springer, 2024, pp. 65–79.
 - [22] Q. Lou, W.-j. Lu, C. Hong, L. Jiang, Falcon: Fast spectral inference on encrypted data, *Advances in Neural Information Processing Systems* 33 (2020) 2364–2374.
 - [23] C. Juvekar, V. Vaikuntanathan, A. Chandrakasan, {GAZELLE}: A low latency framework for secure neural network inference, in: *27th USENIX security symposium (USENIX security 18)*, 2018, pp. 1651–1669.
 - [24] H. Zheng, M. Ji, H. Wang, Y. Liu, L. Fang, Crossnet: An end-to-end reference-based super resolution network using cross-scale warping, in: *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 88–104.
 - [25] D. Kim, C. Guyot, Optimized privacy-preserving cnn inference with fully homomorphic encryption, *IEEE Transactions on Information Forensics and Security* 18 (2023) 2175–2187.
 - [26] J. H. Ju, J. Park, J. Kim, M. Kang, D. Kim, J. H. Cheon, J. H. Ahn, Neujeans: Private neural network inference with joint optimization of convolution and fhe bootstrapping, in: *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security*, 2024, pp. 4361–4375.
 - [27] A. Ebel, K. Garimella, B. Reagen, Orion: A fully homomorphic encryption framework for deep learning, in: *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, Volume 2, 2025, pp. 734–749.
 - [28] R. Dathathri, O. Saarikivi, H. Chen, K. Laine, K. Lauter, S. Maleki, M. Musuvathi, T. Mytkowicz, Chet: an optimizing compiler for fully-homomorphic neural-network inferencing, in: *Proceedings of the 40th ACM SIGPLAN conference on programming language design and implementation*, 2019, pp. 142–156.
 - [29] B. Li, D. Micciancio, On the security of homomorphic encryption on approximate numbers, in: *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, Springer, 2021, pp. 648–677.
 - [30] B. Li, D. Micciancio, M. Schultz-Wu, J. Sorrell, Securing approximate homomorphic encryption using differential privacy, in: *Annual International Cryptology Conference*, Springer, 2022, pp. 560–589.
 - [31] J.-W. Lee, H. Kang, Y. Lee, W. Choi, J. Eom, M. Deryabin, E. Lee, J. Lee, D. Yoo, Y.-S. Kim, et al., Privacy-preserving machine learning with fully homomorphic encryption for deep neural network, *IEEE Access* 10 (2022) 30039–30054.
 - [32] H. Chen, I. Chillotti, Y. Song, Improved bootstrapping for approximate homomorphic encryption, in: *Annual International Conference on the Theory and Applications of Cryptographic Techniques*, Springer, 2019, pp. 34–54.
 - [33] E. Lee, J.-W. Lee, J.-S. No, Y.-S. Kim, Minimax approximation of sign function by composite polynomial for homomorphic comparison, *IEEE Transactions on Dependable and Secure Computing* 19 (2021) 3711–3727.
 - [34] X. Jiang, M. Kim, K. Lauter, Y. Song, Secure outsourced matrix computation and application to neural networks, in: *Proceedings of the 2018 ACM SIGSAC conference on computer and communications security*, 2018, pp. 1209–1222.
 - [35] S. Wang, H. Huang, Secure outsourced computation of multiple matrix multiplication based on fully homomorphic encryption, *KSII Transactions on Internet and Information Systems (TIIS)* 13 (2019) 5616–5630.
 - [36] W.-j. Lu, S. Kawasaki, J. Sakuma, Using fully homomorphic encryption for statistical analysis of categorical, ordinal and numerical data, *Cryptology ePrint Archive* (2016).
 - [37] H. Huang, H. Zong, Secure matrix multiplication based on fully homomorphic encryption, *The Journal of Supercomputing* 79 (2023) 5064–5085.
 - [38] A. Krizhevsky, G. Hinton, Learning multiple layers of features from tiny images, *Technical Report UTML TR 2009*, University of Toronto, Toronto, ON, Canada, 2009.
 - [39] J. H. Cheon, K. Han, A. Kim, M. Kim, Y. Song, Bootstrapping for approximate homomorphic encryption, in: *Advances in Cryptology—EUROCRYPT 2018: 37th Annual International Conference on the Theory and Applications of Cryptographic Techniques*, Tel Aviv, Israel, April 29-May 3, 2018 Proceedings, Part I 37, Springer, 2018, pp. 360–384.
 - [40] Lattigo v5, Online: <https://github.com/tuneinsight/lattigo>, 2023. EPFL-LDS, Tune Insight SA.
 - [41] J. Park, M. J. Kim, W. Jung, J. H. Ahn, Aespa: Accuracy preserving low-degree polynomial activation for fast private inference, *arXiv preprint arXiv:2201.06699* (2022).



Byeong-Seo Min Byeong-Seo Min received the B.S. degree from the School of Computer Science and Engineering, Chung-Ang University, Seoul, South Korea, in 2025. He is currently pursuing the M.S. degree in the Department of Electrical Engineering at Pohang University of Science and Technology (POSTECH), Pohang, South Korea. His research interests include privacy-preserving machine learning and homomorphic encryption.



Joon-Woo Lee Joon-Woo Lee received the B.S. and Ph.D. degrees in Electrical and Computer Engineering from Seoul National University, South Korea, in 2016 and 2022, respectively. He is currently an assistant professor in the Department of Computer Science and Engineering at Chung-Ang University, South Korea. His research interests include privacy-preserving machine learning, homomorphic encryption, post-quantum cryptography, and multiparty computation.