# FastFHE: Packing-Scalable and Depthwise-Separable CNN Inference Over FHE

Wenbo Song[1,3]  Xinxin Fan[1,2]  Quanliang Jing[1]  Shaoye Luo[1,2]  Wenqi Wei[3]  Chi Lin[4]  Yunfeng Lu[5]  Ling Liu[6]
[1]Institute of Computing Technology, CAS    [2]UCAS    [3]Fordham University    [4]Dalian University of Technology
[5]Beihang University    [6]Georgia Institute of Technology
songwenbo@mail.dlut.edu.cn {fanxinxin,jingquanliang, luoshaoye}@ict.ac.cn
c.lin@dlut.edu.cn wenqiwei@fordham.edu lyf@buaa.edu.cn lingliu@cc.gatech.edu

*Abstract*—The deep learning (DL) has been penetrating daily life in many domains, how to keep the DL model inference secure and sample privacy in an encrypted environment has become an urgent and increasingly important issue for various security-critical applications. To date, several approaches have been proposed based on the Residue Number System variant of the Cheon-Kim-Kim-Song (RNS-CKKS) scheme. However, they all suffer from high latency, which severely limits the applications in real-world tasks. Currently, the research on encrypted inference in deep CNNs confronts three main bottlenecks: i) the time and storage costs of convolution calculation; ii) the time overhead of huge bootstrapping operations; and iii) the consumption of circuit multiplication depth.

Towards these three challenges, we in this paper propose an efficient and effective mechanism FastFHE to accelerate the model inference while simultaneously retaining high inference accuracy over fully homomorphic encryption. Concretely, our work elaborates four unique novelties. First, we propose a new scalable ciphertext data-packing scheme to save the time and storage consumptions. Second, we work out a depthwise-separable convolution fashion to degrade the computation load of convolution calculation. Third, we figure out a BN dot-product fusion matrix to merge the ciphertext convolutional layer with the batch-normalization layer without incurring extra multiplicative depth. Last but not least, we adopt the low-degree Legendre polynomial to approximate the nonlinear smooth activation function SiLU under the guarantee of tiny accuracy error before and after encrypted inference. Finally, we execute multi-facet experiments to verify the efficiency and effectiveness of our proposed approach, and the experimental results show that, at the standard 128-bit security level, our FastFHE enables the popular ResNet20 to achieve a 2.41× reduction in inference latency and a 2.38× descent in amortized runtime (using 30 threads) compared to the state-of-the-art methods. Moreover, the performance on several other representative deep CNN architectures, such as ResNet32, ResNet44 and VGG11, also exhibits remarkable advantages of our work in both inference overhead and accuracy.

## 1. Introduction

In today's digital era, spurred by the rapid advancement of artificial intelligence in recent years, machine learning has achieved significant breakthroughs in both software and hardware. As a result, it has gradually become the future direction of development across a wide range of industries. As data volume grows and hardware computing power continues to arise, Machine Learning as a Service (MLaaS) has become a preferred choice for many professionals [1]. MLaaS is a cloud-based service whose primary goal is to lower the entry barriers to machine learning, enabling users from various fields to focus on their business demand without worrying about hardware resources, algorithm implementation, or deployment details [2], [3]. However, MLaaS not only encounters conventional adversarial attacks [4]–[6], but it also suffers from privacy-specific breaches both on sample data and model information, such as membership inference attack [7]–[10], model inversion [11], and model extraction [12]–[14]. One main concern of the clients is the sample data leakage, that is, uploading sensitive data to the cloud server may bring in huge risks of information leakage and privacy breach, which profoundly prompts researchers to look into the issue of privacy-preserving machine learning (PPML) in real-world applications.

PPML is a technical paradigm designed to accomplish machine learning tasks while ensuring data privacy. It addresses the issue of sensitive-data leakage in the pipeline of machine learning, such as medical records, financial transactions, personal behavioral data, etc. Surrounding the core PPML techniques, the fully homomorphic encryption (FHE), which keeps the data encrypted throughout the entire computation process, provides the highest level of privacy protection in theory. It has increasingly become the choice for many researchers to execute PPML.

As the achievement of deep neural networks, secure inference for deep-learning models in a FHE environment is becoming a promising privacy-preserving MLaaS solution at present. As we know, FHE supports performing arbitrary homomorphic additions and homomorphic multiplications on ciphertext without the risk of decryption errors [15], [16]. Fig. 1 briefly illustrates the process of secure deep CNN inference over FHE. That is, a user encrypts private data using the FHE public key and sends the encrypted data to the server in the cloud. The server performs deep CNN inference on the encrypted data and returns the encrypted inference result to the user/client, who would decrypt it with the owned private key to obtain the final plaintext outcome.

Throughout this pipeline, the user's private data remains encrypted during the entire inference process, ensuring the server side cannot access any data information.
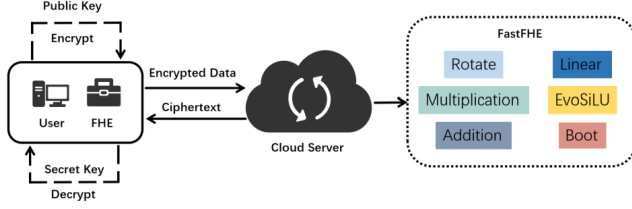


Figure 1: MLaaS over FHE. The dashed box highlights the homomorphic operations in our FastFHE.

Currently, scholars have predominantly focused on constructing deep CNNs over FHE, for example, Gilad-Bachrach et al. [17] and Badawi [18] had successfully deployed neural networks in the encrypted environment. Nevertheless, the practicality is severely constrained by the shallow network depth and the absence of non-linear activation functions, this is due to the available multiplicative depth is limited in leveled homomorphic encryption (Leveled-HE). Subsequently, the FHE is employed to build deeper neural networks, and the bootstrapping operation [16] becomes necessary, this is because it can refresh both the noise level and the multiplicative depth of ciphertext computation, which enables to effectively transform the Leveled-HE into FHE and make it possible to construct arbitrarily deep neural networks. In the reference [19], the bootstrapping operation was used for the first time to realize a deep neural networks over FHE. However, the approach employed a high-degree polynomial to approximate the activation function ReLU, which consumed considerable multiplicative depth and required huge bootstrapping operations, ultimately resulting in a prolonged inference time.

As well known, the information propagation undergoes a long path along layers in the deep CNNs, thus, the operating in encrypted data over FHE accordingly costs a substantial time, a regular and valid solution is to optimize the procedure of convolution computation via an appropriate data-packing scheme, with the purpose of reducing the overall overhead. For instance, GAZELLE [20] offers an efficient FHE-matched convolution algorithm that significantly lowers the number of homomorphic operations needed for traditional convolution in the ciphertext domain. However, in deeper neural networks, the increasing number of input/output channels still leads to high overhead. Furtherly, building on GAZELLE [20], Lee [21] and Lee [19] propose more efficient convolution algorithms and succeed in implementing deeper neural networks. Nonetheless, due to the high computation cost in convolution and the substantial multiplicative depth required for the activation function ReLU, their approaches still demand more overhead.

In addition, for the commonly-used nonlinear activation function, e.g. ReLU, given the FHE only supports homomorphic addition and multiplication operations, hence, the approximation computation on activation functions is usually to use high-degree polynomial. Nevertheless, this causes another problem, i.e., such high-degree polynomial consumes a large amount of multiplicative depth, leading to more frequent bootstrapping operations and significantly prolonging overhead. For instance, Lee et al. [19] used degree-5/15/27 polynomials to approximate the activation function, resulting in a large number of bootstrapping operations. AESPA [22] addresses this challenge by proposing a low-degree polynomial to reduce the overhead of approximately computing ReLU. Despite the overhead is reduced, the inference accuracy declines dramatically. Similarly, Ao and Boddeti [23] adopts a multi-degree polynomial associated with the layered structure, although the optimization on the placement of bootstrapping enables to minimize the amount of bootstrapping operations during convolution, it cannot benefit both efficiency and inference accuracy, still requiring considerable time cost.

To address the highly time-consuming problems while simultaneously retaining the inference accuracy, we propose a ciphertext data packing-scalable and block-ciphertext depthwise-separable CNN inference mechanism FastFHE. Briefly, our work focuses on the core-contents of PPML, i.e. data-packing scheme, multi-channel ciphertext convolution, and approximation expression of activation function. Our solution can achieve a fast and accurate inference for deep CNNs over FHE. Four main contributions are involved:

- **Block-Ciphertext Depthwise-Separable Multi-Channel Convolution.** Traditional convolution under ciphertext grows increasingly time-consuming as the number of output channels rises. Towards this problem, our work proposes a novel scalable data-packing scheme to remarkably reduce the convolution overhead by approximately a factor of $2/(K^2+1)$ ($K$ is the kernel size) through optimizing the usage of rotation keys and rearranging the accumulation of feature values. The more the channels, the more pronounced the reduction in overall cost.

- **Low-Degree Legendre Polynomial Approximation of Activation Function.** In PPML environment, the currently existing methods usually approximate the commonly-used activation function ReLU using either high-degree Chebyshev polynomial, which requires excessive multiplicative depth, or low-degree Hermite polynomial, which compromises model inference accuracy. We think it is the ReLU's sharp corner (non-smooth) that makes it difficult to achieve good approximation with low-degree polynomial. Towards this problem, we adopt a low-degree Legendre polynomial to approximate other smooth activation function-SiLU, to gain a balance between acceptable multiplicative depth consumption and highly-accurate approximation.

- **Fusing Ciphertext Convolutional Layer with Batch-Normalization Layer.** Without increasing extra multiplicative depth consumption, we propose a new architecture ConvBN to fuse the ciphertext convolutional layer and batch-normalization layer,

which not only reduces the overhead of linear-layer multiplicative depth to a large extent, but also notably descends the overhead of activation-layer multiplicative depth by low-degree polynomial.

- **Multi-Facet Experiment Validations on Efficiency and Inference Accuracy.** We implement extensive experiments on ResNet20/32/44 [24] and VGG11 [25] using RNS-CKKS scheme [26], [27]. For ResNet20, the results show the inference latency is 763s with an amortized execution time of 26s in a 30-thread environment, and simultaneously the accuracy remains high after encrypted inference. Compared to the state-of-the-art MPCNN [21], AESPA [22], AutoFHE [23], our FastFHE reduces the latency by 3.26×, 2.51×, 2.41×, and lowers the amortized time by 3.19×, 2.46×, 2.38×, respectively.

## 2. Background and Problem Statement

### 2.1. Homomorphic-Operation Consumption

FHE is an ideal lattice-based encryption mechanism and characterized by the ability to perform computation on encrypted data [15], the core of which lies in allowing to perform addition and multiplication operations on ciphertext, achieving data privacy-preserving. Given the advantage that FHE enables to continuously execute computation on encrypted data until the decryption is required, it has gained more and more attention in the field of PPML. RNS-CKKS [26], [27] is recently proposed to compute floating-point data in FHE field. It integrates CKKS encryption scheme with Residue Number System (RNS), and significantly improves the efficiency of parallel operations. As well known, the CKKS is a homomorphic encryption scheme designed for floating-point computation, it encodes multiple data elements into a polynomial form within a single ciphertext and supports homomorphic addition and multiplication. Meanwhile, through splitting data into multiple parts, RNS effectively reduces the complexity of homomorphic operations.

In CKKS, a plaintext $m$ is first mapped to a polynomial in the ring $\mathbb{Z}_q[X]/\left(X^D + 1\right)$. Then, it is combined with the public key $p_k = (a, e)$ to perform addition and multiplication operations, producing the ciphertext $c = a \cdot m + e$. Here, $D$ is the ring dimension, $a$ and $e$ are noise terms from the public key. The original polynomial is further split across several prime moduli, by which the parallel computation can be realized. The partial results from multi-modulus computation are merged to obtain the final outcome, in this way, the overall computation efficiency can be improved to a large extent. The decryption is achieved through key-switching procedure, and the ciphertext noise from the individual prime modulus is transformed into the original domain. For each two ciphertexts $ct_1 = (c_{10}, c_{11})$ and $ct_2 = (c_{20}, c_{21})$, the homomorphic addition and multiplication are defined as:

$$ct_{add} = (c_{10} + c_{20}, c_{11} + c_{21}) \bmod q, \qquad (1)$$

$$ct_{mult} = (d_0 + d_1) + RLK(d_2) \bmod q, \qquad (2)$$

TABLE 1: Execution time for various homomorphic operations, wherein "Pt" and "Ct" denote plaintext-ciphertext and ciphertext-ciphertext, "CRot" denotes ciphertext rotation, "Add" and "Mul" indicate homomorphic addition and multiplication, "Boot" represents bootstrapping.

| Operation | AddPt | AddCt | MulPt | MulCt | CRot | Boot |
|-----------|-------|-------|-------|-------|------|------|
| Time(s) | 0.20s | 0.22s | 0.57s | 1.26s | 1.06s | 14.176s |

where $d_0 = (c_{10} \cdot c_{20}) \bmod q$, $d_1 = (c_{10} \cdot c_{21} + c_{11} \cdot c_{20}) \bmod q$, $d_2 = (c_{11} \cdot c_{21}) \bmod q$, $RLK(d_2)$ represents the relinearization operation, which applies the relinearization key $rlk$ to transform $d_2$ into a linear form.

Another important concept in FHE is the multiplicative depth closely tied to the ring dimension $D$, i.e. the larger the dimension $D$, the greater the maximum multiplicative depth $L$. Each homomorphic multiplication consumes one level of multiplicative depth. Once all levels are used up, a bootstrapping operation [16] is required to refresh the ciphertext's multiplicative depth, its validity and precision. Table 1 shows our experimental results on the time costs of various homomorphic operations under RNS-CKKS scheme, it reveals the bootstrapping is significantly more time-consuming than the other homomorphic operations. Therefore, reducing bootstrapping as much as possible is critical for efficient encrypted computation.

**Tremendous Bootstrapping and Numerous Rotation operations Are Required in Convolution Calculation.** During the convolution calculation over FHE, the data arrangement in ciphertext slots is equivalently important, as it affects the number of rotation and bootstrapping operations. The rotation refers to cyclically shifting the encrypted data for the task execution, such as rearranging data, performing encryption computation using specific patterns. On the other hand, when many homomorphic multiplications are needed, each ciphertext may require bootstrapping to continuously implement encrypted-data computation, resulting in frequent bootstrapping operations. Therefore, reducing the number of ciphertexts is beneficial to the efficiency. Furthermore, if a single ciphertext contains a large amount of data, the complex calculation would necessitate numerous rotations, as shown in Table 1, it is also time-consuming as well.

To sum up, upon the statement above, we know that, in order to improve efficiency, it is crucial to strike a balance between packing large amount of data into fewer ciphertexts and minimizing costly rotation and bootstrapping operations.

### 2.2. Existing Ciphertext-Convolution Solution

For image-specific deep CNNs, such as object recognition, the convolution calculation plays a pivotal role in extracting local feature from the input data. For example, a 2D image is first converted into a three-dimensional tensor and then passed through convolutional filters/kernels. Let $\{c_i, h_i, w_i\}$ and $\{c_o, h_o, w_o\}$ represent the input and output dimensions of a convolution layer, where $c$ indicates the number of channels and $h$, $w$ represent the height and width.
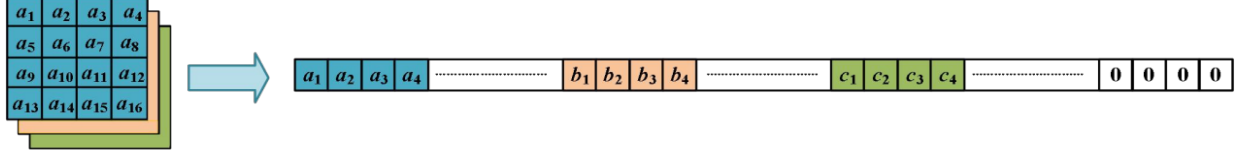
Figure 2: Packing image into ciphertext slots, wherein $a_i$, $b_i$, and $c_i$ denote feature values respectively from three channels. Any unused store slots are filled with zeros.

The filter size in convolution is denoted as $(f \times f)$. Given ResNet [24] as a classic CNN architecture to achieve a deep learning in the true sense for the first time by introducing "deep residual connections", we subsequently discuss our work referring to ResNet as an example, and set a stride of $s = 1, 2$ and padding to 1 in the convolutional layers.

Over FHE, the data in ciphertext slots is stored as one-dimensional vector. Hence, we must flatten the three-dimensional tensor data of the input image into a one-dimensional vector. Then, we can fill these ciphertext slots based on a specific data arrangement strategy as shown in Fig. 2. Such a ciphertext data-packing method can straightly influence the computation efficiency of convolutional layer. As we know, the maximum size of feature map in ResNet20 is $16 \times 32 \times 32 = 16384$, thus, the number of store slots on-demand is set to $2^{14} = 16384$.

Towards the complex convolution calculation, Gazelle [20] proposed an efficient scheme SISO for the deep CNN inference over FHE, in which the $h_i \times w_i$ feature values are flattened and packed into a single ciphertext, along with the slot arrangement as shown in Fig. 2.

**Stationary Data-Packing and Idle Slots Raise Computation Complexity.** Investigating on the currently popular SISO scheme, we find the following observations and problems: i) since the input data is packed in a single ciphertext, the only way to execute convolution is to rotate the ciphertext. As a result, both the input ciphertext and convolution kernel need to be rotated according to different rotation indices, by which $f_h \times f_w$ rotated ciphertexts and kernel plaintexts are generated. The numbers of rotations and specific indices depend on two factors: the size of defined kernels and whether the feature values participate in the convolution operation; ii) each rotated ciphertext is then multiplied by its concordant rotated kernel plaintext, and the products from the $f_h \times f_w$ rotations are accumulated to produce the SISO convolution result; iii) take CIFAR-10 dataset and ResNet [24] as an example, each image sample has 3 input channels (RGB), and the channel count gradually increases. Using SISO for the convolution operation would produce a large number of ciphertexts, leading to huge bootstrapping operations, which significantly extends the total computation time of deep CNNs over FHE; and iv) on the other hand, when a large stride is set, e.g. $s = 2$, the ciphertext obtained after convolution would contain more idle slots, which causes more influence on the efficiency.

From these observations, we unveil two aspects that affect the computational complexity, i.e. the stationary block-packing scheme and idle slots incurred by larger stride.

### 2.3. Approximation of Activation Function

Given that the ciphertext computation in FHE only supports homomorphic addition and multiplication, hence, we need elaborate an combination of such two operations to approximately compute the nonlinear activation function.

**Accuracy Degradation with Low-Degree Polynomial Approximation.** To achieve the approximation of nonlinear activation function, CryptoDL [28], as one of the foundational works, employs polynomial to represent activation function in a homomorphic encryption (HE) environment. Since the bootstrapping is unavailable under HE, it cannot apply into deeper neural networks. Subsequently, some other works [22], [29]–[32] explore the low-degree polynomial for the activation-function approximation. One of which, AESPA [22] utilizes a degree-2 Hermite polynomial to approximate ReLU and gains a remarkable reduction on runtime. Nevertheless, the common limitation of such category of low-degree polynomial approximation methods lies in they cannot accurately represent the activation function, resulting in significant degradation on inference accuracy.

**Time-Consuming Computation with High-Degree Polynomial Approximation.** Take into consideration the accuracy, on the high-degree polynomial side, Lee et al. [19] compare diverse configurations and find a combination of degree-5/15/27 minimax approximation polynomials can obtain an appropriate accuracy. This allows pre-trained models to be directly used, nonetheless, it still consumes substantial multiplicative depth stemming from numerous bootstrapping operations within the residual blocks. Towards this expensive cost problem, AutoFHE [23] recently frames this problem as a multi-objective optimization, and employs a layered and mixed-degree polynomial to seek balance between maximum accuracy and minimized number of bootstrapping operations, nevertheless, the large and sophisticated composition still demands extremely-high computational resources, e.g., at least 768GB of RAM, which makes it hard to deploy in regular environment without high-performance computing resources.

Hence, a question emerges naturally, i.e. whether there exists a proper low-degree approximation polynomial to achieve a remarkable overhead reduction, and simultaneously maintain high inference accuracy.

### 2.4. Threat Model

As stated in AESPA [22], MPCNN [21] and AutoFHE [23], our proposed solution is applicable to their mentioned

threat model as well. Specifically, we address the PPML for users' data in MLaaS scenarios, that is, when users rely on cloud server to execute CNN inference, they anticipate their sensitive data to be secure and not disclosed to the server. To this end, they usually encrypt their private data using a public key and upload both the ciphertext and public key to the server. The server then performs encrypted-data computation and inference on the ciphertext, and returns the computation result to the users. Then, they decrypt with their private key to obtain the final inference outcome. Throughout the entire process, the server only deals the ciphertext and cannot access to the private (plaintext) data.

## 3. FastFHE: An Efficient and Effective PPML

### 3.1. Feature Map-Aligned Block-Oriented Packing

Diverse homomorphic operations have different computational overheads as reported in Table 1, from which we see the bootstrapping costs the most time. Thereby, we need to minimize the frequency of bootstrapping, and a direct and effective solution is to reduce the number of ciphertexts, especially when multiple channels exist in each layer. As we know, the SISO [20] performs convolution calculation on a single-input/output channel, instead, our work aims to extend to multi-channel convolution and further reduce the cost of ciphertext rotation and multiplication operations. In detail, our approach employs a scalable $N\times N$ block-based data-packing scheme regarding the variational sizes of feature maps during the pipeline of convolution, by which all channels of feature maps can be packed into a single ciphertext. Of course, the ciphertext capacity, i.e. the number of slots, must accommodate the size of the largest feature map in the whole convolution pipeline. For example, upon the analysis on ResNet20, we know the largest feature map is: $16\times32\times32 = 16384$. As well known, for a CNN, the semantics-extracting process is usually to narrow down the size of feature map and extend the number of channels. Routinely, the stride may not be set too large, this is because large stride would cause much loss of local information, thus, the downsampling is usually to scale feature map to the half of original size, and correspondingly double the channel number. To comply with 2D matrix block, we defines the data-packing block parameter $N$ as:

$$N = \begin{cases} \dfrac{(F_{max})^{\frac{1}{2}}}{h_i} & s = 1 \\[3mm] \dfrac{(F_{max})^{\frac{1}{2}}}{h_o} & s = 2 \end{cases}, \qquad (3)$$

where $F_{max}$ denotes the size of the largest feature map, e.g. 128 for ResNet20, $h_i$ and $h_o$ are input channel's height and output channel's height, and $s$ is the stride of convolution.

We can comparatively analyze our method with the prevalent SISO scheme [20], each position that stores a single element in SISO is currently treated as an $N\times N$ block. The size of block per row depends on the stride, and each $N\times N$ block stores the feature values from all channels.



(a) $3\times32\times32$ data packing

(b) $16\times32\times32$ data packing

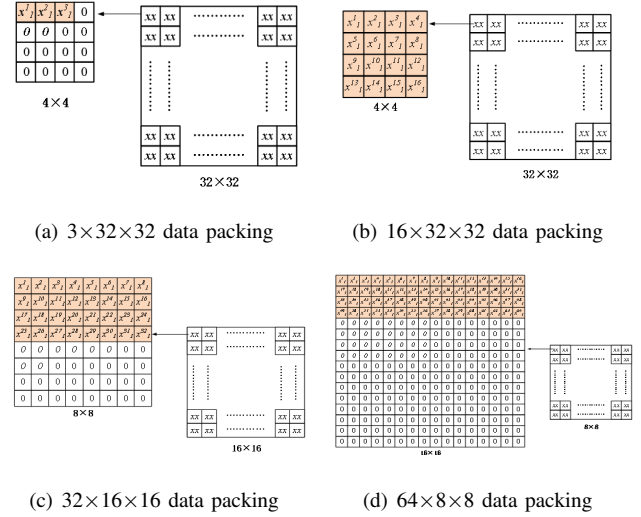(c) $32\times16\times16$ data packing

(d) $64\times8\times8$ data packing

Figure 3: Scalable block-oriented packing scheme.

Once the number of channels is less than block size, the rest elements of the bock are padded with zeros as placeholders, keeping a complete square. Also take ResNet20 as an example, as shown in Fig. 3, the ciphertext has $2^{14}=16384$ slots, which can be viewed as a $128\times128$ square. This square is composed of multiple $N\times N$ blocks. For a $16\times32\times32$ feature map, each $4\times4$ block can store the 16-channel information at a single position. Similarly, Fig. 3(b)-(d) give the continued downsampling layer by layer from feature-map size $32\times32$ with channel number 16 to $16\times16$ with channel number 32, and to $8\times8$ with channel number 64. In this way, the multi-channel data is neatly packed into a single ciphertext. That is to say, our proposed multi-channel convolution can be carried out by manipulating just one ciphertext, which significantly reduces the computation and storage complexities, compared to the stationary SISO-packing scheme. Next, we design homomorphic convolution fashion to match our new data-packing scheme.

### 3.2. Ciphertext Depthwise-Separable Convolution

**Traditional Convolution Manipulation.** As stated previously, for each feature map during convolution, our method need have the capability to tackle $N\times N$ block ciphertext as input, and produce another $N\times N$ block ciphertext as output. Prior to clarifying our proposed block-ciphertext depthwise-separable convolution, we at first introduce the routine operation on ciphertext. Taking a $4\times4$ feature map and $3\times3$ convolution kernel as an example, we demonstrates the detailed steps for the traditional convolution manipulation in Fig. 11 in Appendix A.

The traditional convolution scans the input data with a set of convolution kernels to extract local features. For a deep CNN, the beginning convolutional layers can capture shallow features, such as edges, textures and colors, while the deep layers enable to extract the complex and abstract

feature. Therefore, in such a deep CNN, as the number of kernels or input channels grows, the computation load will become enormous, or even the bottleneck over FHE.

**Block-Ciphertext Depthwise-Separable Convolution Manipulation.** The traditional convolution slides the kernel over the input image, implementing an elementwise multiplication and summation at each position to produce the output feature map. If there are multiple input channels and convolutional kernels, each kernel executes convolution with its associated channel. Then, the results are summed to form the final output feature map as shown in Fig. 4(a). According to such a routine procedure to perform convolution, the computation load would be extremely expense.

To address this issue, as shown in Fig. 4(b), a depthwise-separable convolution fashion is figured out in consideration of FHE, that is, the traditional convolution is decomposed into two smaller parts: depthwise convolution and pointwise convolution, by which both the parameter amount and computation load are remarkably cut down. Depthwise convolution acts as the first stage of our proposed depthwise-separable convolution, wherein each input channel is convolved independently with a dedicated spatial kernel. Unlike traditional convolution, it employs the same number of kernels as input channels, with each kernel operating solely on its associated channel and no cross-channel fusion occurring. This design significantly reduces computation load, as the convolution of each channel is processed independently while preserving the original channel count. Pointwise convolution constitutes the second stage, utilizing $1\times1$ kernel to perform cross-channel feature fusion. Since the kernel size is $1\times1$, it retains the spatial dimension of the feature map while only modifying the channel depth, whose primary role is to linearly combine the multi-channel outputs from depthwise convolution, projecting them into a new channel space to enable efficient inter-channel communication with minimal computation load.

**Analysis on Parameter Amount and Computation Load.** We formally analyze the saving amounts of parameters and computation load. Suppose the kernel size is $f$, the numbers of input and output channels are $c_i$ and $c_o$, the output feature map has a size of $h_o\times w_o$. For the traditional convolution, the number of parameters is $f\times f\times c_i\times c_o$, and the total computation amount is $f\times f\times c_i\times c_o\times h_o\times w_o$. However, for our proposed depthwise-separable convolution, the total number of parameters is $f\times f\times c_i + c_i\times c_o$, and the total computation amount becomes $f\times f\times c_i\times h_o\times w_o + c_i\times c_o\times h_o\times w_o$. Accordingly, the ratio of parameters between the proposed depthwise-separable convolution and the traditional convolution is:

$$\frac{f\times f\times c_i + c_i\times c_o}{f\times f\times c_i\times c_o} = \frac{1}{c_o} + \frac{1}{f^2}, \qquad (4)$$

and the ratio of computation load between them is:

$$\frac{f\times f\times c_i\times h_o\times w_o + c_i\times c_o\times h_o\times w_o}{f\times f\times c_i\times c_o\times h_o\times w_o} = \frac{1}{c_o} + \frac{1}{f^2}. \qquad (5)$$

When the output channel $c_o$ is relatively large (e.g. 32 or 64), the term $\frac{1}{c_o}$ can be ignored. If $f = 3$ (i.e. a

$3\times3$ convolution kernel), under our depthwise-separable convolution, the computational load drops to around $\frac{1}{9}$ of that of the traditional convolution. Next, we will analyze how to apply in an encrypted environment.

**Quantifying Ciphertext Computation in Traditional Convolution.** In the ciphertext domain, the traditional convolution process is illustrated in Fig. 11 in Appendix A. At first, both the encrypted input data and the encoded kernel plaintext must be rotated, which requires $(f^2 - 1)$ rotations. Then, the corresponding ciphertext and plaintext are multiplied and summed to compute the convolution results across multiple input channels, which involves $f^2$ multiplications and consuming 1-level multiplicative depth. Afterward, these multi-channel convolution results are summed to generate the output feature map for a single output channel, which requires $c_i$ rotations and 1 multiplication, and consumes one more multiplicative depth. Thus, for one output channel, the total rotation cost is $f^2 - 1 + c_i$, the total multiplication cost is $f^2 + 1$, and the total levels of multiplicative depth cost is 2. To produce $c_o$ output channels, the required operations become $(f^2 - 1 + c_i)\times c_o$ rotations, $(f^2+1)\times c_o$ multiplications and 2-level multiplicative depth. Fig. 5 sketches the procedure of block-ciphertext depthwise-separable multi-channel convolution, and the more concrete description is illuminated in Appendix B.

**Quantifying Ciphertext Computation in Depthwise-Separable Convolution.** In our depthwise-separable convolution under ciphertext scenario, the encrypted input data and encoded kernel plaintext at first undergo a depthwise convolution, which requires $(f^2 - 1)$ rotations and $f^2$ multiplications. Afterward, the pointwise convolution is performed, with the cost of $c_i$ rotations and 2 multiplications. For $c_o$ output channels, the process in total uses $(f^2 - 1 + c_i\times c_o)$ rotations and $(f^2 + 2c_o)$ multiplications. Therefore, in the ciphertext environment, the ratio of required rotations between our proposed depthwise-separable convolution and the traditional convolution is:

$$\frac{\text{Depthwise-Separable Convolution}}{\text{Traditional Convolution}} = \frac{f^2 - 1 + c_i\times c_o}{(f^2 - 1 + c_i)\times c_o}, \qquad (6)$$

and the ratio of multiplications between them is:

$$\frac{f^2 + 2c_o}{(f^2 + 1)\times c_o} = \frac{f^2}{(f^2 + 1)\times c_o} + \frac{2}{f^2 + 1} \approx \frac{2}{f^2 + 1}. \qquad (7)$$

It is clear that as the number of output channels increases, the depthwise-separable convolution achieves an increasingly pronounced reduction on computation load. Given that the beginning conventional layers capture the shallow features of the input samples, such as edges, textures, and colors, which play important roles for the downstream tasks, hence, we apply such block-ciphertext depthwise-separable multi-channel convolution fashion to those layers that have a large number of channels. When $c_o$ is large, the $\frac{f^2}{(f^2+1)\times c_o}$ term can be neglected, the number of multiplications is cut down to $\frac{1}{5}$ ($f = 3$) of the original approximately. Furthermore, regrading the parameters of kernel plaintext, for instance, a $32\times16\times16$ tensor would
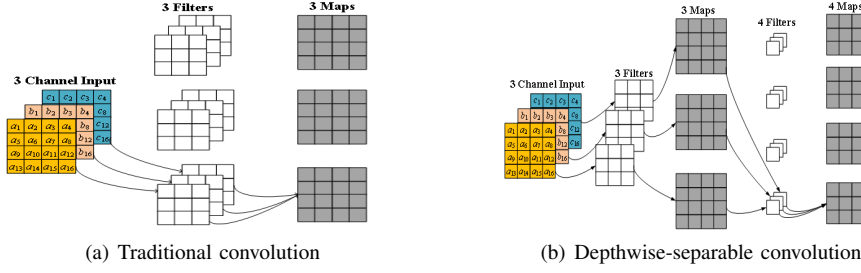
(a) Traditional convolution  (b) Depthwise-separable convolution

Figure 4: Architectures of traditional convolution vs. depthwise-separable convolution.



(a) Depthwise convolution



(b) Pointwise convolution



(c) Accumulation of the input-channel convolution results to form a single output-channel feature map



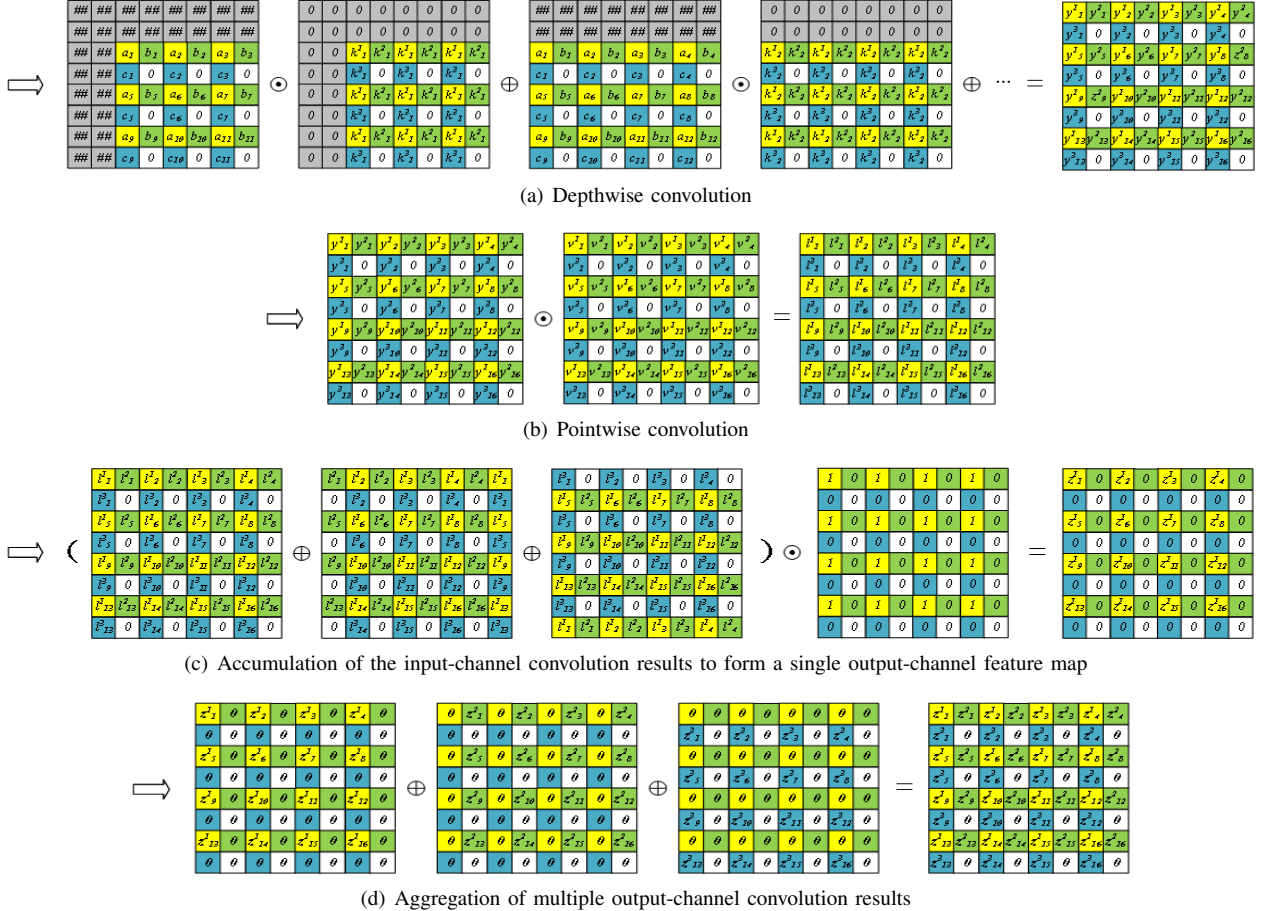(d) Aggregation of multiple output-channel convolution results

Figure 5: Block-ciphertext depthwise-separable multi-channel convolution.

require $f^2 \times c_o = 9 \times 32 = 288$ plaintexts of convolution kernel (parameter files) in the traditional convolution, but only $f^2 + c_i = 9 + 32 = 41$ plaintexts in our proposed depthwise-separable convolution.

## 3.3. Fusing Convolutional Layer with BN Layer

As shown in Fig. 5, splitting the single-step convolution process into two parts of (depthwise and pointwise) convolutions entails one additional multiplicative depth. In FHE, the utilization of multiplicative depth must be care-fully planned to control both the placement and frequency of bootstrapping with the purpose of keeping the added overhead within a reasonable range. Taking into account that the convolutional layer and batch-normalization (BN) layer are both the combination of multiplication and addition, one question appears naturally, whether we can merge them together to reduce multiplications and multiplicative depths or not. To this end, we propose a BN dot-product fusion matrix with the purpose of fully taking advantage of the homomorphic multiplication and associated multiplicative depth in the convolution layers. That is to say, we can merge

the convolutional layer and BN layer while simultaneously eliminating the extra multiplicative depth incurred by the proposed depthwise-separable convolution operations.

To accomplish the fusion, we first reshape the BN layer's calculation formula into a form that consists of multiplication and addition operations only, yielding:

$$y_i = \gamma \hat{x}_i + \beta = \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$
$$= \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}(x_i - \mu_B) + \beta = \alpha \mathbb{X}_i + \beta, \qquad (8)$$

where $\alpha = \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}$ represents the BN layer's coefficient, and $\mathbb{X}_i = x_i - \mu_B$ is the result after subtracting the batch mean from the convolution output. In this way, we obtain a BN layer formula that is expressed solely in terms of addition and multiplication. Specifically, $x_i$ is first subtracted by the mean $\mu_B$, then multiplied by the coefficient $\alpha$, and finally the bias $\beta$ is added.

**ConvBN without Multiplicative-Depth Optimization.** One necessary manipulation is to integrate the BN layer's computation into the convolutional layer's calculation process. As demonstrated in Fig. 5, after the execution of pointwise convolution, the multi-input channel convolution results are accumulated to form a single output-channel result, and then a mask is applied to filter out the redundant data produced by ciphertext rotations. This process consumes two levels of multiplicative depth. Our approach aims to merge these steps into one and incorporate the BN layer, forming a ConvBN layer. Appendix C illustrates the ConvBN process without the multiplicative-depth optimization. Compared with Fig. 5, the only modification lies in that the result obtained from the pointwise convolution is first subtracted by the mean matrix $\mu_B$, then multiplied by the BN coefficient $\alpha$, which combines the mask matrix with the BN layer's coefficient matrix, finally, during the aggregation of multi-output channel results, the offset matrix $\beta$ is added.

**ConvBN with Multiplicative-Depth Optimization.** On the basis of Fig. 12 in Appendix C, we can further optimize the multiplicative depth so that only one level is consumed. In detail, our approach is to fuse the parameters of pointwise convolution with the BN coefficient $\alpha$ to form a new coefficient matrix—the BN dot-product fusion matrix. Then, through performing the ciphertext rotations, multiplications and accumulations between the depthwise convolution and BN dot-product fusion matrix, we can simultaneously accumulate the convolution results from different input channels and apply BN manipulation for the regularization purpose, the procedure is described in Fig. 6. Therefore, our proposed BN dot-product fusion matrix achieves the aggregation of pointwise convolution, BN coefficient $\alpha$ and mask to filter out redundant data.

Consequently, the element of BN dot-product fusion matrix is the product of the pointwise convolution kernel parameter and the BN coefficient, namely,

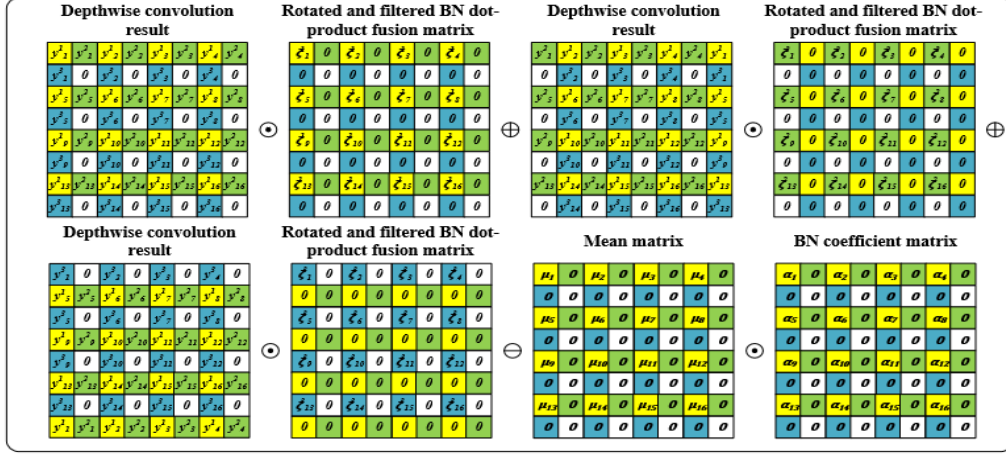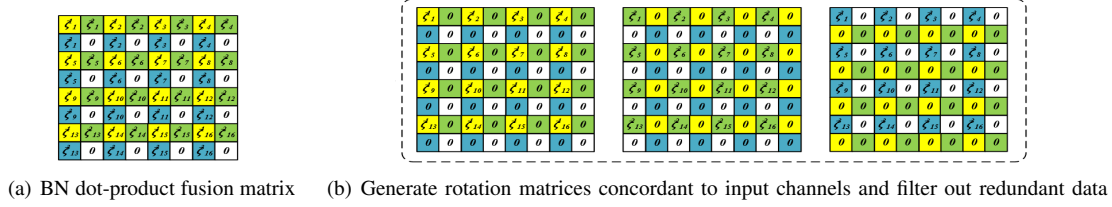$$\zeta_n^i = \nu_n^i \times \alpha_o = \nu_n^i \times \frac{\gamma}{\sqrt{\sigma_B^2 + \epsilon}}, \qquad (9)$$

where $\alpha_o$ denotes the BN coefficient for each output channel's BN layer, and $\nu_n^i$ represents the pointwise convolution kernel parameter for the $i$th input channel at position $n$. After generating the BN dot-product fusion matrix, we need to perform both rotation and filtering on this matrix. The rotation, performed according to the number of input channels, serves to shift the parameters of each channel to the first position within the $N \times N$ block. Consequently, when the depthwise convolution result undergoes the corresponding rotation, multiplying it by the rotated fusion matrix directly produces the result of the pointwise dot-product, followed by the BN coefficient multiplication for that channel, and the result is placed in the first position within the block to facilitate the accumulation of multi-input channel convolution results. Furthermore, the filtering step is designed to eliminate any redundant data brought by the rotations of depthwise convolution. We provide the schematic diagrams of our proposed ConvBN with/without optimization on multiplicative depth consumption in Appendix D.

Since the parameters of BN dot-product fusion matrix are derived from plaintext, this matrix does not incur any limitation in terms of multiplicative depth consumption or rotation count, that is to say, it can be generated directly in the plaintext domain. The actual consumption of multiplicative depth and rotations within the ciphertext is confined to the content inside the solid-line box in Fig. 6. This portion consolidates the two steps in Fig. 5(b)-(c) of the block-ciphertext depthwise-separable convolution into a single step by utilizing rotated and filtered BN dot-product fusion matrix. As a result, under the same rotation count, one multiplicative depth operation can be saved.

## 4. Activation Function Approximation

Since FHE only handles the homomorphic multiplication and addition operations, scholars usually leverage polynomial to approximate the nonlinear activation functions. At present, the pioneering CryptoDL [28] and other subsequent works [30]–[33] have employed different-degree polynomials to approximate the activation functions, however, few of them have demonstrated the effectiveness (inference accuracy) and rationality in consideration of both efficiency and approximation precision. Upon the in-depth investigation on existing studies, we find the challenge lies in the following predicament: the low-degree polynomial cannot precisely approximate the activation function, further leading to a significant drop in model inference accuracy; while the high-degree polynomial consumes a huge amout of multiplicative depth, enforcing frequent bootstrapping operations in residual blocks in spite that a high inference accuracy can be gained. Therefore, an ideal solution is to seek a lightweight solution that enables to precisely approximate the nonlinear activation function without consuming too much multiplicative depth at the same time.

At present, the high-degree polynomial typically employs Chebyshev polynomial to approximate ReLU function, while the low-degree polynomial usually utilizes

(a) BN dot-product fusion matrix

(b) Generate rotation matrices concordant to input channels and filter out redundant data

(c) Fusion of steps (b) and (c) in Fig. 5 with optimized multiplicative depth consumption

Figure 6: BN dot-product fusion matrix and ConvBN multiplicative depth optimization.

Hermite polynomial. As well known, the polynomial-approximation technique can perform well on smooth functions, whereas the ReLU function has a discontinuity point, it is not a smooth function. To gain an accurate approximation, the Chebyshev polynomial must be high-degree. Oppositely, the Hermite polynomial with low-degree inevitably sacrifices inference accuracy. To address the dilemma, we replace ReLU with SiLU, a continuous function over real number domain $\mathbb{R}$ and mathematically expressed as,

$$SiLU(x) = x \cdot \sigma(x) = \frac{x}{1 + e^{-x}}, \quad (10)$$

where $\sigma(x) = \frac{1}{1+e^{-x}}$ denotes the standard Sigmoid function. Since it provides gradient over its entire range, especially in the negative region, it prevents the neurons from completely ceasing to update. Consequently, it tends to perform better than ReLU and its variants (Leaky ReLU).

Aiming to realize lightweight computation, we adopt low-degree Legendre polynomials to approximate SiLU function, stemming from the property of smoothness. The goal of Legendre polynomial approximation is to minimize the mean squared error, that is, to minimize the integral of the squared error over the entire range:

$$err = \int_{-1}^{1} (f(x) - \sum_{k=0}^{n} a_k P_k(x))^2 dx, \quad (11)$$

where $f(x)$ denotes activation function, and $\sum_{k=0}^{n} a_k P_k(x)$ indicates the polynomial for approximation. Thus, it pro-

vides a better average approximation over the entire interval, and this excellent performance across a wide range is clearly more suitable for representing nonlinear activation functions.

Legendre polynomials are a class of orthogonal polynomials that play a crucial role in mathematics and physics, typically denoted as $P_n(x)$, where $n$ represents the degree of the polynomial. These polynomials can be computed using the following recurrence relation:

$$P_n(x) = \begin{cases} 1 & n = 0 \\ x & n = 1 \\ \frac{(2n+1)xP_n(x) - nP_{n-1}(x)}{n+1} & n \geq 2 \end{cases}. \quad (12)$$

For a given function $f(x)$, we can express it as a linear combination of Legendre polynomials:

$$f(x) = \sum_{n=0}^{\infty} a_n P_n(x), \quad (13)$$

where $a_n$ is determined by the following formula:

$$a_n = \frac{2n+1}{2} \int_{-1}^{1} f(x) P_n(x) dx. \quad (14)$$

By computing sufficiently many coefficients $a_n$, we can obtain an approximation polynomial with a finite number of terms. Since our goal is to achieve a precise approximation and consume minimal multiplicative depth, we used the formula above to derive a polynomial of maximum degree 5 to approximate SiLU. This approximation consumes 3 levels of multiplicative depth, as shown in Fig. 7.
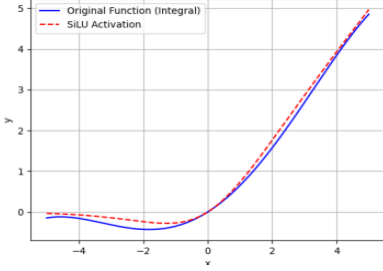
Figure 7: Polynomial approximation of SiLU

## 5. Reshaping Network Architecture over FHE

By integrating our proposed ConvBN and employed activation function SiLU, and subsequently incorporating downsampling, AVGPOOL, fully connected layers, we need reshape the network architecture using RNS-CKKS scheme, and take ResNet20 as an example as shown in Fig. 8.

Since the traditional convolution can fully extract the shallow features and our ConvBn is especially effective to reduce computation load when the number of channels is large, we use conventional convolution in Layer 1 and adopt the proposed ConvBn in Layer 2 and Layer 3. We determine the placement of bootstrapping based on two factors: i) whether the remaining multiplicative depth is insufficient to support the following operation, otherwise the bootstrapping must be performed; and ii) ensuring the ciphertexts involved in skip connections are at the same multiplicative depth as much as possible, otherwise the addition of ciphertexts at different multiplicative depths would incur computation error or decryption failure.

## 6. Performance Evaluation

### 6.1. Configuration

To validate the efficiency and effectiveness of our proposed FastFHE, we perform multi-facet experiments on the commonly-used dataset CIFAR-10 [34] using RNS-CKKS scheme from OpenFHE Library [35]. The detailed configurations on the running environment, dataset and models, settings on encryption-parameters, statistics of three baselines and our FastFHE, are demonstrated in Appendix E.

### 6.2. Resource Consumption Performance

For the encrypted-data computation, the consumption of various resources is an important evaluation criterion, we assess the performance in terms of ciphertext size, runtime memory usage, and key size. As shown in Table 2, the pioneering work [19] leverages single-ciphertext and single-channel (SCSC) data-packing scheme to apply into PPML domain, while AESPA [22] and AutoFHE [23] use single-ciphertext and multi-channel (SCMC) data-packing manner. Easy to know, the SCSC data-packing scheme dominates more computing resources compared to the SCMC, this is because the ciphertext amount will be enlarged as the number of channels increases, for example, the former consumes 376MB to store ciphertext and 473 GB memory to run procedure, while the latter consumes much smaller storage spaces, i.e. 8.5 MB for AESPA and 7.2MB for AutoFHE, and uses smaller memories, i.e. 316 GB and 270GB respectively. Compared to these baselines, our FastFHE consumes the least resources, i.e. 6.3MB for ciphertext storage and 80GB memory for procedure running. Moreover, our approach significantly reduces the key size to less than $\frac{1}{2}$ of that of the best baseline.

Furthermore, AESPA [22] and AutoFHE [23] mainly focus the attention on the approximation of activation layers, ignoring the pivotal correlation between data arrangement and rotation manipulation. This oversight would induce excessively high memory consumption, deriving from the huge number of ciphertext rotations to calculate convolution in the pipeline of neural network, in addition to the loading of rotation key and the homomorphic multiplication operations. To address this issue, our work proposes the scalable ciphertext data-packing scheme, it enables the multi-channel encrypted data to be arranged within each block in a particular rule, by which both the number of ciphertext rotations and the amount of rotation keys can be remarkably reduced.

In addition, our proposed data-packing scheme can also achieve the simultaneous shifting of multiple channels in a single ciphertext rotation, greatly declining the memory overhead. As shown in Fig. 9, compared to the baselines, our FastFHE needs the smallest amount of rotation keys and memory usage under ResNet20. On the other hand, during the ciphertext convolution operations, the rotated ciphertexts produced by homomorphic rotations also occupy additional memory space, the more rotations performed, the larger the memory footprint of these rotated ciphertexts. These two factors are the fundamental causes of high memory consumption observed in the three baseline models.

### 6.3. Overhead in Depthwise-Separable Convolution

This section compares our proposed block-ciphertext depthwise-separable convolution (BCDSC) with the traditional ciphertext convolution. In the early stages of inference using ResNet20, the input features are abstract, and the neural network needs to fully capture both spatial and inter-channel characteristics to uncover correlations between different channels. Traditional convolution is more effective at extracting inter-channel information. Therefore, ciphertext convolution is employed at the initial layer and at Layer 1. In the later stages, where the number of channels increases significantly, the proposed block-ciphertext depthwise-separable convolution is adopted to reduce computational complexity and parameter amount. This treatment aims to improve convolution efficiency in ciphertext domain while maintaining inference accuracy at the same time.

Under the ciphertext environment, we process CIFAR-10 test data using our scalable $N \times N$ block-oriented data-packing method. The proposed BCDSC is applied at Layer
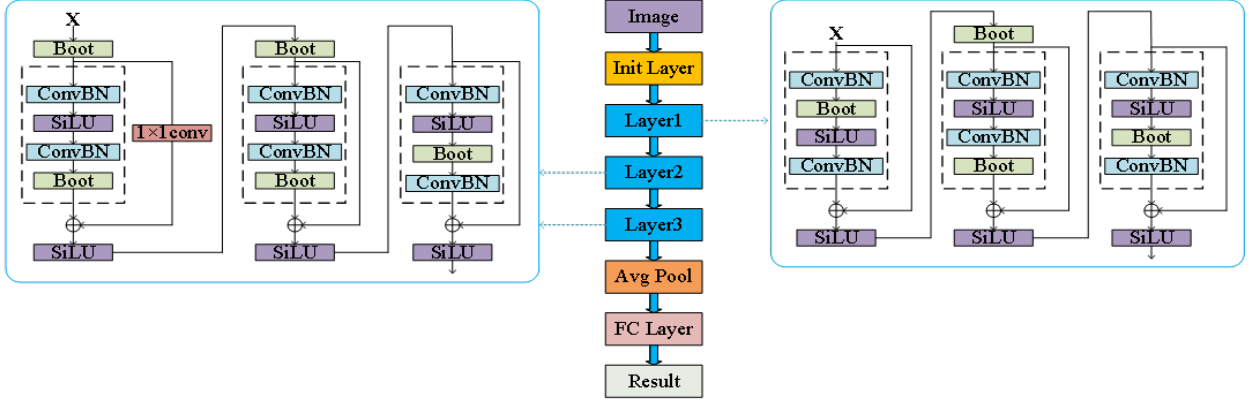
Figure 8: Rebuilt ResNet20 architecture in consideration of FHE.

TABLE 2: Resource consumption with ResNet20.

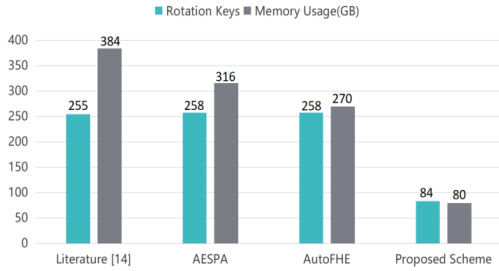| Model | Packing Method | Ciphertext Size | Memory Usage | Key Size |
|---|---|---|---|---|
| Literature [19] | Single-ciphertext, single-channel | 376MB | 472GB | 109GB |
| AESPA [22] | Single-ciphertext, multi-channel | 8.5MB | 316GB | 133GB |
| AutoFHE [23] | Single-ciphertext, multi-channel | 7.2MB | 270GB | 121GB |
| FastFHE | $N \times N$ Block based packing | **6.3**MB | **80**GB | **42.1**GB |



Figure 9: Rotation keys and memory usage.

2 and Layer 3, and the activation functions are approximated by a degree-6 Chebyshev polynomial resembling the ReLU function, ensuring compatibility with ciphertext-domain CNN models. Under both traditional convolution and our BCDSC, Table 3 and Table 4 detail the statistics for each layer and total time for convolution operations, as well as the number of plaintext convolution kernel parameters involved. Since the number of parameters directly correlates with the number of multiplications, it also indirectly reflects the computation advantage of our proposed method.

During traditional convolution, the multi-channel ciphertexts packed in $N \times N$ blocks need to generate $f^2$ rotated ciphertexts. These rotated ciphertexts are multiplied and accumulated with the corresponding $N \times N$ block packed multi-channel plaintext kernels to compute the convolution results across input and output channels. Therefore, for one convolutional layer, the traditional method requires $f^2 \times c_o$ plaintexts of convolution kernels. Specifically, Layer 2 requires 1782 plaintexts, and Layer 3 requires 3456 plaintexts. In contrast, under the proposed BCDSC scheme, deep

convolution only requires $f^2$ plaintexts, while pointwise convolution requires $c_o$ plaintexts. The total plaintext in Layer 2 is $(f^2 + c_o) \times n_{conv} = 246$, and in Layer 3 is 438. As shown in Table 4, since the initial layer and Layer 1 still use traditional convolution under ciphertext, the number of plaintexts required in those stages is the same. However, in Layer 2 and Layer 3, the proposed method significantly reduces the number of plaintexts, which also implies a substantial reduction in the number of homomorphic multiplication operations.

As previously analyzed in Equ. 7, the number of multiplications required by our scheme is approximately $\frac{2}{f^2+1} = \frac{1}{5}$ of that by the traditional convolution. In Table 4, the plaintext amount under our scheme in Layer 2 and Layer 3 is reduced by approximately 5 times in theory compared to traditional convolution (by $6.52\times$ and $7.01\times$ experimentally, due to the different number of external multiplications required for matrix-vector multiplication in the calculation). Furthermore, as shown in Table 3, the execution time for our proposed BCDSC in Layer 2 and Layer 3 is only $\frac{1}{4.86}$ and $\frac{1}{5.03}$ of that of traditional ciphertext convolution. This fully demonstrates the advantages of our proposed block-ciphertext depthwise-separable multi-channel convolution fashion in terms of inference efficiency.

## 6.4. Model Inference Latency

Upon the optimization on ConvBN multiplicative depth and low-degree Legendre approximation function, we employ ResNet20, ResNet32, ResNet44, and VGG11 [25] to execute our proposed FastFHE. Table 5 presents the inference latency and amortized runtime, wherein the fastest variant of AutoFHE [23] is used in our experiments.

TABLE 3: Convolution calculation time.

| Convolution Algorithm | Init. Layer | Layer 1 | Layer 2 | Layer3 | Total |
|---|---|---|---|---|---|
| Traditional convolution method | 8s | 54s | 472s | 754s | 1,327s |
| Block ciphertext depthwise separable convolution method | 8s | 54s | 97s | 150s | 309s |

TABLE 4: Number of plaintext convolution kernel parameters.

| Convolution Algorithm | Init. Layer | Layer 1 | Layer 2 | Layer 3 | Total |
|---|---|---|---|---|---|
| Traditional convolution method | 144 | 864 | 1782+32 | 3,456+64 | 6,342 |
| Block ciphertext depthwise separable convolution method | 144 | 864 | 246+32 | 438+64 | 1,788 |

Note: "+32" and "+64" at Layer 2 and Layer 3 denote the plaintexts for downsampling convolutions in skip connections

Seen from the results, we have the observations: i) it is evident MPCNN exhibits longer inference latency and amortized runtime compared to AESPA and AutoFHE, this is primarily because MPCNN does not pay specific attention on activation function approximation. It employs a high-degree (6, 7 or 8) Chebyshev polynomial function to approximate ReLU function; ii) in contrast, AESPA employs the low-degree (2) Hermite polynomial to approximate ReLU function. AutoFHE leverages layered mixed-degree polynomial approximations for the substitution of activation functions. It is the low/mixed-degree polynomial that decreases the multiplicative depth, then it reduces the frequency of bootstrapping operations. AutoFHE has similar inference latency and amortized runtime with AESPA, however, whose accuracy loss is less than AESPA as shown in Table 7, stemming from the variable polynomial's degree settings; and iii) although AESPA and AutoFHE focus on lightweight treatment on the activation layers, they ignore the long-time inference problem in linear layers. Differently, through adopting the low-degree Legendre polynomials to approximate the employed SiLU in our work, FastFHE incurs lower time consumption in the activation layer, furtherly, coupled with the efficient convolution based on scalable data-packing scheme, our FastFHE significantly outperforms the baselines in both inference latency and amortized runtime. Specifically, in a 30-thread environment, our amortized runtime is reduced by 3.19×, 2.46×, and 2.38× compared to MPCNN, AESPA, and AutoFHE, respectively.

**Runtime of various Operation.** This subsection is to experimentally evaluate the execution time of various operations on the four reshaped neural network architectures. Table 6 shows the time overhead for various operations, from which the following concludes can be obtained: i) for MPCNN, the activation function is approximated by high-degree polynomial, leading to much higher consumption compared to HerPN and EvoReLU. Moreover, to achieve high approximation precision, more frequent bootstrapping operations are required (accounting for 71% to 74% of the total inference time), which is the main reason for the high inference latency; ii) for AESPA, the time consumption for HerPN operations is lower. However, this low-degree polynomial results in remarkably reduced model inference accuracy as shown in Table 7, a trade-off unacceptable in exchange for shorter inference times, especially in the

security-critical scenario (e.g. self-driving); iii) for the comparison purposes, we select the fastest variant of AutoFHE, where it exhibits linear layer latency similar to AESPA, but differs in the activation function, this is because AutoFHE employs a layered mixed-degree polynomial, and consumes slightly more time than AESPA, yet delivers a less accuracy loss compared to AESPA; iv) our FastFHE utilizes the depthwise-separable convolution, which drastically reduces the computation load of convolution operations, in addition to ConvBN optimization, as a result, its linear-layer time consumption is notably lower than that of other schemes. Additionally, by approximating SiLU function with low-degree Legendre polynomial, our FastFHE not only achieves an excellent approximation to maintain high inference accuracy, but also reduces the multiplicative depth consumption. This contributes significantly to lowering the time required for both activation functions and bootstrapping operations, thereby, further reducing the overall inference latency and amortized runtime.
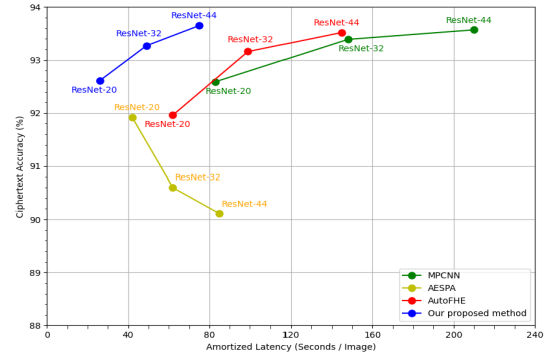


Figure 10: Cipertext accuracy and amortized latency.

## 6.5. Model Inference Accuracy

On the basis of comprehensive utilization of our proposed scalable data-packing scheme, ConvBN optimization, and low-degree Legendre polynomial approximation, we run a group of experiments to observe the overall comparison in terms of ciphertext accuracy and amortized latency on CIFAR-10 dataset under ResNet20/32/44. Obviously, our FastFHE located in the top-left corner in Fig. 10 has the best performance compared to the three baselines.

TABLE 5: Amortized runtime per image under multiple threads (30 threads).

| Backbone | | MPCNN | | AESPA | | AutoFHE | | FastFHE | |
|---|---|---|---|---|---|---|---|---|---|
| Model | Parameters | Latency | Amortized Latency | Latency | Amortized Latency | Latency | Amortized Latency | Latency | Amortized Latency |
| ResNet-20 | 269K | 2487s | 83s | 1913s | 64s | 1842s | 62s | **763**s | **26**s |
| ResNet-32 | 464K | 4438s | 148s | 3068s | 101s | 2944s | 99s | **1463**s | **49**s |
| ResNet-44 | 658K | 6302s | 210s | 4775s | 147s | 4346s | 145s | **2012**s | **75**s |
| VGG-11 | 123K | 1369s | 46s | 572s | 19s | 656s | 22s | **272**s | **9**s |

TABLE 6: Runtime of various operations.

| Backbone | MPCNN | | | AESPA | | | AutoFHE | | | FastFHE | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Linear | AppReLU | Boot | Linear | HerPN | Boot | Linear | EvoReLU | Boot | Linear | EvoSiLU | Boot |
| ResNet-20 | 417s | 279s | 1764s | 899s | **23**s | 870s | 863s | 49s | 897s | **350**s | 35s | **383**s |
| ResNet-32 | 852s | 448s | 2929s | 1539s | 53s | 1476s | 1519s | 91s | 1333s | **647**s | **45**s | 747s |
| ResNet-44 | 906s | 623s | 4629s | 2503s | **82**s | 2490s | 2472s | 117s | 1756s | **804**s | 106s | 1102s |
| VGG-11 | 177s | 167s | 998s | 347s | 8s | 226s | 457s | 6s | 174s | **114**s | **5**s | 153s |

TABLE 7: Accuracy and error between plaintext and ciphertext inference.

| Backbone | | MPCNN | | AESPA | | AutoFHE | | FastFHE | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Model | Accuracy | Accuracy | Error | Accuracy | Error | Accuracy | Error | Plaintext | Ciphertext | Error |
| ResNet20 | 92.76% | 92.59% | 0.17% | 91.92% | 0.84% | 91.96% | 0.80% | 92.66% | 92.61% | **0.05%** |
| ResNet32 | 93.51% | 93.39% | 0.12% | 90.60% | 2.91% | 93.16% | 0.35% | 93.34% | 93.27% | **0.07%** |
| ResNet44 | 93.92% | 93.57% | 0.35% | 90.11% | 3.81% | 93.52% | 0.40% | 93.86% | 93.65% | **0.21%** |
| VGG11 | 90.53% | 90.41% | **0.12%** | 88.97% | 1.56% | 90.69% | 0.16% | 90.22% | 90.01% | 0.21% |

Although our proposed depthwise-separable convolutions in Layer 2 and Layer 3 can result in somewhat loss of inter-channel feature extraction, the efficient low-degree polynomial approximation of activation function, as well as the superior performance of smooth activation function SiLU under plaintext inference, both can compensate the accuracy. Moreover, due to the high-quality approximation of the activation function, the gap between plaintext accuracy and ciphertext accuracy is narrowed down to a tiny extent. As shown in Table 7, when using ResNet20 for inference, the accuracy error of our FastFHE before and after encryption is only 0.05, whereas the accuracy errors of MPCNN, AESPA, and AutoFHE are 0.17%, 0.84%, and 0.80%, respectively. Similarly, under ResNet32 and ResNet44, our method also yields the smallest accuracy loss, i.e., 0.07% and 0.21%. For VGG11, the error is slightly higher than that of MPCNN and AutoFHE, this is primarily because ReLU outperforms our employed SiLU in the shallow layers.

## 7. Related Work

We review the existing work from two perspectives: i) the computation-load reduction of linear layer; and ii) the approximate expression of activation function of nonlinear layer. The detailed discussions and comparative analytics are presented in Appendix F.

## 8. Conclusion

We have experimentally and analytically validated the efficiency and effectiveness of our proposed FastFHE, a privacy-preserving deep neural network mechanism, through figuring out a scalable data-packing scheme, a lightweight yet highly-precise low-degree approximation polynomial, and a ConvBN architecture by fusing convolutional layer with BN layer. To conquer the grave drawbacks of tremendous bootstrapping and rotation operations, high resource consumption caused by the currently-prevalent stationary data-packing scheme, as well as the dilemma of time consumption and model inference accuracy in present works, we propose a novel scalable data-packing scheme to save the time consumption and storage resources, and employ depthwise-separable convolution fashion to reduce the computing load of convolution calculation. Furthermore, we design BN dot-product fusion matrix to integrate the convolution layer with the BN layer without incurring extra multiplicative depth. In addition, we identify a fundamental challenge in polynomial approximation on nonlinear activation functions lies in the difficulty of precisely approximating a non-smooth function like ReLU. Towards this issue, we adopted low-degree Legendre polynomial to approximately express the smooth function SiLU. Finally, extensive experiments demonstrate that, compared to the baselines MPCNN, AESPA and AutoFHE, our proposed FastFHE achieves the inference latency reductions of $3.26\times$, $2.51\times$, and $2.41\times$ udner ResNet20, and in a 30-thread environment, the amortized runtime is reduced by $3.19\times$, $2.46\times$ and $2.38\times$, respectively. From the additional experiments on ResNet32, ResNet44 and VGG11, the analogous analytics can be gained as well. Thus, our proposed lightweight PPML solution can boost the realistic applications in near future.

# References

[1] J. Huang, Y. Xu, Q. Wang, Q. C. Wang, X. Liang, F. Wang, Z. Zhang, W. Wei, B. Zhang, L. Huang, J. Chang, L. Ma, T. Ma, Y. Liang, J. Zhang, J. Guo, X. Jiang, X. Fan, Z. An, T. Li, X. Li, Z. Shao, T. Qian, T. Sun, B. Diao, C. Yang, C. Yu, Y. Wu, M. Li, H. Zhang, Y. Zeng, Z. Zhang, Z. Zhu, Y. Lv, A. Li, X. Chen, B. An, W. Xiao, C. Bai, Y. Mao, Z. Yin, S. Gui, W. Su, Y. Zhu, J. Gao, X. He, Y. Li, G. Jin, X. Ao, X. Zhai, H. Tan, L. Yun, H. Shi, J. Li, C. Fan, K. Huang, E. Harrison, V. C. Leung, S. Qiu, Y. Dong, X. Zheng, G. Wang, Y. Zheng, Y. Wang, J. Guo, L. Wang, X. Cheng, Y. Wang, S. Yang, M. Fu, and A. Fei, "Foundation models and intelligent decision-making: Progress, challenges, and perspectives," *The Innovation*, vol. 6, no. 6, p. 100948, 2025.

[2] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.* OpenReview.net, 2021.

[3] A. Fawzi, M. Balog, A. Huang, T. Hubert, B. Romera-Paredes, M. Barekatain, A. Novikov, F. J. R. Ruiz, J. Schrittwieser, G. Swirszcz *et al.*, "Discovering faster matrix multiplication algorithms with reinforcement learning," *Nature*, vol. 610, no. 7930, pp. 47–53, 2022.

[4] B. Wang, X. Fan, Q. Jing, H. Tan, and J. Bi, "Advcgan: An elastic and covert adversarial examples generating framework," in *International Joint Conference on Neural Networks, IJCNN 2021, Shenzhen China, July 18-22, 2021.* IEEE, 2021, pp. 1–8.

[5] B. Wang, X. Fan, Q. Jing, Y. Su, J. Li, and J. Bi, "GTAT: adversarial training with generated triplets," in *International Joint Conference on Neural Networks, IJCNN 2022, Padua, Italy, July 18-23, 2022.* IEEE, 2022, pp. 1–8.

[6] S. Luo, X. Fan, Q. Jing, C. Lin, M. Li, Y. Lu, and Y. Xu, "Sifternet: A generalized and model-agnostic trigger purification approach," *10.48550/arXiv.2505.14531*, 2025.

[7] H. Liu, Y. Wu, Z. Yu, and N. Zhang, "Please tell me more: Privacy impact of explainability through the lens of membership inference attack," in *IEEE Symposium on Security and Privacy, SP 2024, San Francisco, CA, USA, May 19-23, 2024.* IEEE, 2024, pp. 4791–4809.

[8] J. Li, N. Li, and R. Bruno, "MIST: defending against membership inference attacks through membership-invariant subspace training," in *33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024*, D. Balzarotti and W. Xu, Eds. USENIX Association, 2024.

[9] Y. He, B. Li, Y. Wang, M. Yang, J. Wang, H. Hu, and X. Zhao, "Is difficulty calibration all we need? towards more practical membership inference attacks," in *Proceedings of the 2024 on ACM SIGSAC Conference on Computer and Communications Security, CCS 2024, Salt Lake City, UT, USA, October 14-18, 2024*, B. Luo, X. Liao, J. Xu, E. Kirda, and D. Lie, Eds. ACM, 2024, pp. 1226–1240.

[10] Z. Chen and K. Pattabiraman, "A method to facilitate membership inference attacks in deep learning models," in *32nd Annual Network and Distributed System Security Symposium, NDSS 2025, San Diego, California, USA, February 24-28, 2025.* The Internet Society, 2025.

[11] S. Shi, N. Wang, Y. Xiao, C. Zhang, Y. Shi, Y. T. Hou, and W. Lou, "Scale-mia: A scalable model inversion attack against secure federated learning via latent space reconstruction," in *32nd Annual Network and Distributed System Security Symposium, NDSS 2025, San Diego, California, USA, February 24-28, 2025.* The Internet Society, 2025.

[12] P. Lv, H. Ma, K. Chen, J. Zhou, S. Zhang, R. Liang, S. Zhu, P. Li, and Y. Zhang, "Mea-defender: A robust watermark against model extraction attack," in *IEEE Symposium on Security and Privacy, SP 2024, San Francisco, CA, USA, May 19-23, 2024.* IEEE, 2024, pp. 2515–2533.

[13] M. Tang, A. Dai, L. DiValentin, A. Ding, A. Hass, N. Z. Gong, Y. Chen, and H. H. Li, "Modelguard: Information-theoretic defense against model extraction attacks," in *33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024*, D. Balzarotti and W. Xu, Eds. USENIX Association, 2024.

[14] T. Nayan, Q. Guo, M. Alduniawi, M. Botacin, A. S. Uluagac, and R. Sun, "Sok: All you need to know about on-device ML model extraction - the gap between research and practice," in *33rd USENIX Security Symposium, USENIX Security 2024, Philadelphia, PA, USA, August 14-16, 2024*, D. Balzarotti and W. Xu, Eds. USENIX Association, 2024.

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

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

[17] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy," in *International conference on machine learning.* PMLR, 2016, pp. 201–210.

[18] A. Al Badawi, C. Jin, J. Lin, C. F. Mun, S. J. Jie, B. H. M. Tan, X. Nan, K. M. M. Aung, and V. R. Chandrasekhar, "Towards the alexnet moment for homomorphic encryption: Hcnn, the first homomorphic cnn on encrypted data with gpus," *IEEE Transactions on Emerging Topics in Computing*, vol. 9, no. 3, pp. 1330–1343, 2020.

[19] 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*, vol. 10, pp. 30 039–30 054, 2022.

[20] C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan, "Gazelle: A low latency framework for secure neural network inference," in *27th USENIX security symposium (USENIX security 18)*, 2018, pp. 1651–1669.

[21] E. Lee, J.-W. Lee, J. Lee, Y.-S. Kim, Y. Kim, J.-S. No, and 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. 12 403–12 422.

[22] J. Park, M. J. Kim, W. Jung, and J. H. Ahn, "Aespa: Accuracy preserving low-degree polynomial activation for fast private inference," *arXiv preprint arXiv:2201.06699*, 2022.

[23] W. Ao and 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.

[24] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016.* IEEE Computer Society, 2016, pp. 770–778.

[25] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015.

[26] J. H. Cheon, K. Han, A. Kim, M. Kim, and Y. Song, "A full rns variant of approximate homomorphic encryption," in *International Conference on Selected Areas in Cryptography.* Springer, 2018, pp. 347–368.

[27] J. H. Cheon, A. Kim, M. Kim, and Y. Song, "Homomorphic encryption for arithmetic of approximate numbers," in *Advances in cryptology–ASIACRYPT 2017: 23rd international conference on the theory and applications of cryptology and information security, Hong kong, China, December 3-7, 2017, proceedings, part i 23.* Springer, 2017, pp. 409–437.

[28] E. Hesamifard, H. Takabi, M. Ghasemi, and R. N. Wright, "Privacy-preserving machine learning as a service," *Proceedings on Privacy Enhancing Technologies*, 2018.

[29] H. Chabanne, A. De Wargny, J. Milgram, C. Morel, and E. Prouff, "Privacy-preserving classification on deep neural network," *Cryptology ePrint Archive*, 2017.

[30] T. Ishiyama, T. Suzuki, and H. Yamana, "Highly accurate cnn inference using approximate activation functions over homomorphic encryption," in *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 3989–3995.

[31] S. Obla, X. Gong, A. Aloufi, P. Hu, and D. Takabi, "Effective activation functions for homomorphic evaluation of deep neural networks," *IEEE access*, vol. 8, pp. 153 098–153 112, 2020.

[32] E. Hesamifard, H. Takabi, and M. Ghasemi, "Deep neural networks classification over encrypted data," in *Proceedings of the Ninth ACM Conference on Data and Application Security and Privacy*, 2019, pp. 97–108.

[33] Q. Lou and L. Jiang, "Hemet: A homomorphic-encryption-friendly privacy-preserving mobile neural network architecture," in *International conference on machine learning*. PMLR, 2021, pp. 7102–7110.

[34] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

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

[36] T. van Elsloo, G. Patrini, and H. Ivey-Law, "Sealion: a framework for neural network inference on encrypted data," *CoRR*, vol. abs/1904.12840, 2019.

[37] L. Rovida and A. Leporati, "Encrypted image classification with low memory footprint using fully homomorphic encryption," *Cryptology ePrint Archive*, 2024.

[38] J. L. Crawford, C. Gentry, S. Halevi, D. Platt, and V. Shoup, "Doing real work with fhe: the case of logistic regression," in *Proceedings of the 6th Workshop on Encrypted Computing & Applied Homomorphic Cryptography*, 2018, pp. 1–12.

[39] A. Brutzkus, R. Gilad-Bachrach, and O. Elisha, "Low latency privacy preserving inference," in *International Conference on Machine Learning*. PMLR, 2019, pp. 812–821.

[40] F. Bourse, M. Minelli, M. Minihold, and P. Paillier, "Fast homomorphic evaluation of deep discretized neural networks," in *Advances in Cryptology–CRYPTO 2018: 38th Annual International Cryptology Conference, Santa Barbara, CA, USA, August 19–23, 2018, Proceedings, Part III 38*. Springer, 2018, pp. 483–512.

# Appendix A.
# Ciphertext Convolution under Traditional Operation

Taking a $4\times4$ feature map and $3\times3$ convolution kernel as an example, as shown in Fig. 11, the following demonstrates the detailed steps for the traditional convolution operation.

- **Step i):** Pack the feature map into pre-defined $N\times N$ blocks to form ciphertexts. Here, we use $2\times2$ blocks for illustration. The multi-channel feature map is then transformed into the block units.
- **Step ii):** Obtain the rotated ciphertexts for the kernel through rotating the ciphertext $f\times w$ times ($f$ and $w$ depend on the convolution kernel size). Multiply and accumulate to get the convolution results for the three input channels, as shown in Fig. 11(a).
- **Step iii):** Combine the convolution results for the input channels to produce the output channel's feature map. The results are packed into $N\times N$ block positions concordant to the rotation offsets. A single reduction step is performed at this point to preserve only valid data through removing any extra rotation-induced artifacts, as illustrated in Fig. 11(b).
- **Step iv):** Repeat Steps i)-iii) for multi-channel settings. The combined results are placed into the associated block positions in ciphertext, generating an $N\times N$ ciphertext containing multi-channel output features, as shown in Fig. 11(c).

# Appendix B.
# $N\times N$ Block Depthwise-Separable Multi-Channel Convolution

As illustrated in Fig. 5, the following presents the detailed procedure of depthwise-separable multi-channel convolution operations

- **Step i):** Multiply the ciphertext and kernel plaintext using the relevant rotation indices, then accumulate the results to obtain depthwise convolution outcome.
- **Step ii):** Perform pointwise convolution to produce the multi-input channel convolution result.
- **Step iii):** Accumulate the results of the multi-input channel convolution via rotations, generating the output feature map for a single output channel.
- **Step iv):** If multiple output channels are required, repeat Steps ii)-iii), placing the generated output into the concordant channel position in the $N\times N$ block. Ultimately, an $N\times N$ block-oriented ciphertext containing multiple output channels can be produced.

# Appendix C.
# ConvBN without Optimization

As shows in Fig. 12, the ConvBN without optimization on multiplicative depth consumption is calculated as follows: the result obtained from the pointwise convolution is first subtracted by the mean matrix $\mu_B$, then multiplied by the BN coefficient matrix $\alpha$, which combines the mask matrix with the BN layer's coefficient matrix, finally, during the aggregation of multi-output channel results, the offset matrix $\beta$ is added.

# Appendix D.
# ConvBN Diagram

Fig. 13 sketches the diagram of our proposed ConvBN with/without optimization on multiplicative depth consumption. For the optimization, the coefficient $\zeta$ denotes the BN dot-product fusion matrix aggregated from the pointwise kernel parameters $\nu$ and the BN coefficients $\alpha$. $\zeta^n$ is the rotated and filtered version of the BN dot-product fusion matrix. After the depthwise convolution result $y$ is rotated, it is multiplied by $\zeta^n$, and the products are accumulated, then, subtract $\mu \times \alpha$ to complete BN layer computation.

# Appendix E.
# Experiment Configuration

**Running Environment.** The server is configured with two Intel Xeon Gold 6430 processors (Base Frequency 2.1 GHz/Max Turbo Frequency 3.4 GHz, 32 cores, 64 threads, 60 MB cache, 270WTDP), 8X 64GB Samsung DDRS ECC REGS 4800 memory modules, Ubuntu 22.04.

**Dataset and Models.** The dataset CIFAR-10 has 10 categories with each comprising 5,000 training images and 1,000 testing images. This dataset is widely used by convolutional neural networks in computer vision tasks, e.g. object recognition, despite the dataset is small-scale, it does not influence our evaluation, this is because our method focuses on the metrics of efficiency and accuracy loss over FHE. Therefore, we can employ the pretrained models in the plaintext environment, such as ResNet20/32/44. The ResNet20's rebuilt architecture is shown in Fig. 8, analogously, the ResNet32/44 can be reshaped as well.

**Encryption-Parameter Settings.** Regarding the encryption parameters towards the reshaped architectures of ResNet20/32/44, we set the ring dimension to $D = 2^{16}$, the number of slots to $n = 2^{14}$, the secret key's Hamming weight to $h = 64$, and the multiplicative depth to 26, with bootstrapping consuming 14 layers of multiplicative depth. The security level is set to $\lambda \geq 128$.

**Baselines.** Since the previous PPML models, such as CryptoNet [17], SEALion [36], CryptoDL [28], were applied into simple machine learning models with small number of layers, or replaced the non-arithmetic activation with simple arithmetic activation disregarding the usage of highly-precise polynomial approximation accompanied with bootstrapping operations. Hence, it is unsuitable to compare with them. We in this paper refer to the most related methods MPCNN [21], AESPA [22] and AutoFHE [23] as the baselines as shown in Table 8, wherein the brief
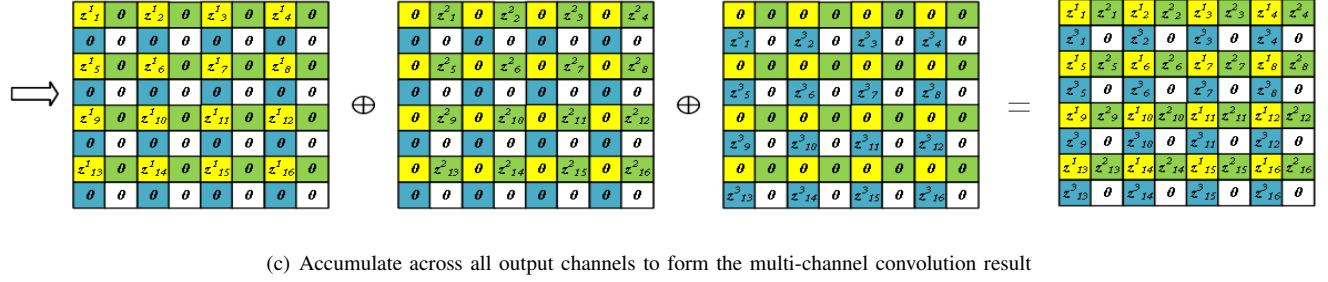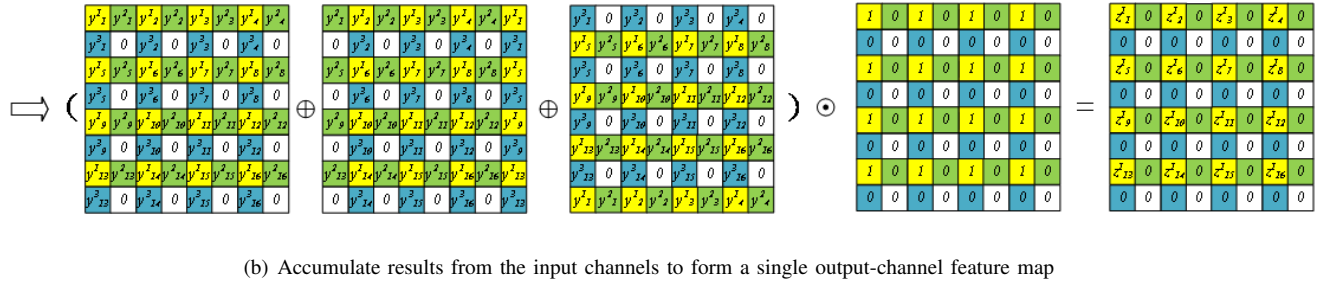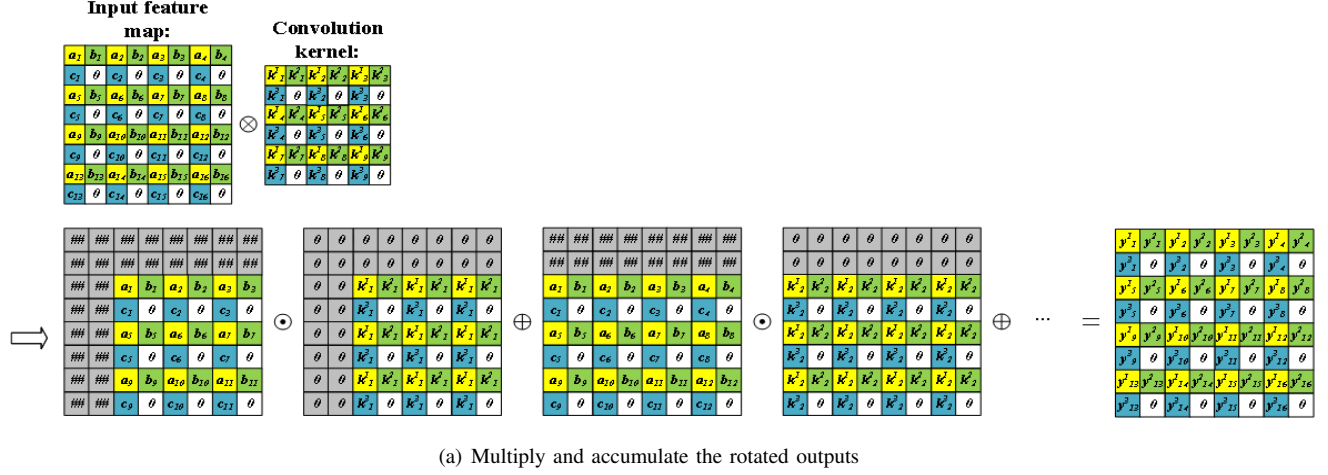
(a) Multiply and accumulate the rotated outputs



(b) Accumulate results from the input channels to form a single output-channel feature map



(c) Accumulate across all output channels to form the multi-channel convolution result

Figure 11: Multi-channel ciphertext convolution under traditional operations.

TABLE 8: Statistics on baselines and our FastFHE.

| Model Name | Venue | Encryption | Polynomial | Activation |
|---|---|---|---|---|
| MPCNN [21] | ICML'22 | CKKS | high-degree | ReLU |
| AESPA [22] | arXiv'22 | CKKS | low-degree | ReLU |
| AutoFHE [23] | USENIX'24 | CKKS | multi-degree | mixed |
| FastFHE | —- | CKKS | low-degree | SiLU |

statistics information is exhibited. To our best knowledge, the related work is limited to date.

MPCNN [21], AESPA [22] and AutoFHE [23] all adopt multi-channel packing schemes along with traditional convolution algorithms. Specifically, MPCNN focuses on the implementation of standard convolution under a single-ciphertext multi-channel setting and leverages high-degree polynomial approximations to improve the expressiveness of the ReLU function, thereby enhancing the inference accuracy over FHE. However, its high computational complexity results in prolonged inference time. AESPA utilizes low-degree Hermite polynomials to approximate the ReLU activation function, remarkably reducing the computational cost over FHE, as well as decreasing the number of bootstrapping operations and overall inference latency, nevertheless, its poor approximation on activation function significantly degrades inference accuracy. AutoFHE designs a selection algorithm that adaptively chooses polynomials of different degrees to approximate the activation function at different network layers. However, it cannot address the high computation cost problem. In their framework, achieving high accuracy still requires a considerable inference time.
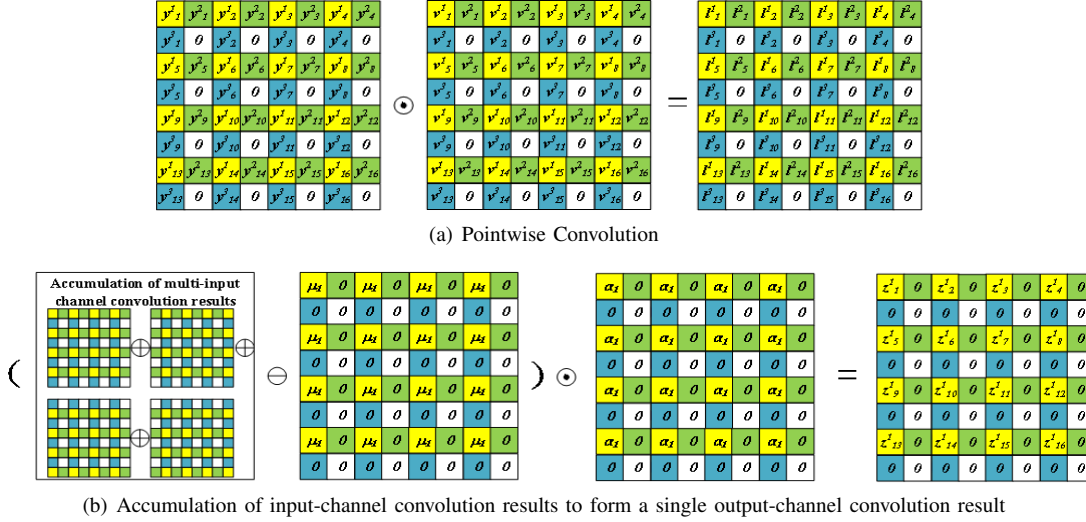
(a) Pointwise Convolution



(b) Accumulation of input-channel convolution results to form a single output-channel convolution result

Figure 12: ConvBN without optimization on multiplicative depth consumption.



(a) ConvBN without multiplicative-depth optimization



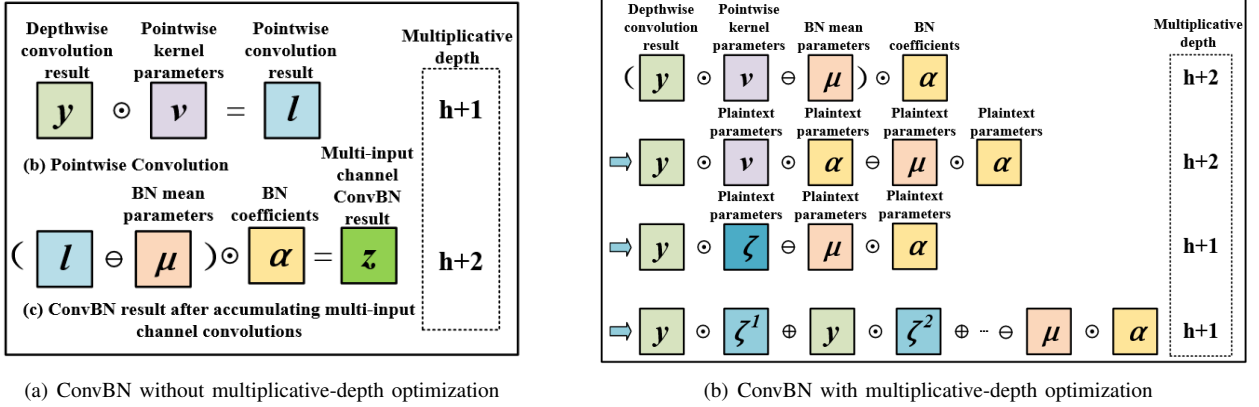(b) ConvBN with multiplicative-depth optimization

Figure 13: Schematic diagram of ConvBN.

# Appendix F.
# Related Work

## F.1. Computation-Load Reduction of Linear Layer

With the advancement of fully homomorphic encryption (FHE) techniques in recent years, Cheon et al. [16] proposed a bootstrapping algorithm under the RNS-CKKS scheme, which made it feasible to perform privacy-preserving machine learning (PPML) over FHE. Since then, researchers have turned their attention to implementing PPML based on RNS-CKKS. However, turning this concept into a practical solution remains challenging. Lee et al. [19] were the first to realize privacy-preserving deep neural network inference using bootstrapping, but their work did not consider the architectural constraints imposed by the FHE environment. As a result, their model suffered from long inference times. To address this, researchers begin focusing on the computational cost of linear layers, which remains one of the primary bottlenecks in FHE-based deep learning inference. Juvekar et al. [20] propose GAZELLE, a framework that combines secure multi-party computation and homomorphic encryption to enable efficient neural network inference. They introduce an effective convolution method under homomorphic encryption that significantly reduces the number of homomorphic operations for standard convolutions. However, in deeper networks, as the numbers of input and output channels grows, the convolution of multiple ciphertexts still incurs high time and memory costs. Building on this, Lee et al. [21] propose a more efficient convolution scheme and successfully implement deep neural networks over FHE. Nonetheless, due to the inherently high computation complexity of traditional convolutions and the need to enlarge encryption parameters for packing multi-channel ciphertext data, their method still requires prolonged inference time and consume substantial memory resources. Rovida et al. [37] introduce an optimized vector-encoded convolution method that requires fewer rotation keys. However, due

to the relatively outdated packing strategies and ciphertext slot utilization, the number of rotation operations increase significantly. Although their method reduces the memory consumption, the excessive homomorphic rotations lead to long computation times for linear layers and negatively impacts the overall model performance.

Despite significant progress has been made in ciphertext-domain convolution, existing approaches still incur substantial computational overhead when performing convolutions on encrypted data. The time-consuming overhead hinders the practical deployment of privacy-preserving inference over FHE. In this regard, our FastFHE addresses the challenge by employing a scalable data-packing scheme for ciphertext, which compactly encodes multi-channel data into a single ciphertext. This strategy significantly reduces the number of homomorphic rotation operations required during linear layer computation, thereby lowering memory consumption. Additionally, our proposed FastFHE introduces a specially designed ConvBN architecture tailored to this packing scheme, further reducing the time cost of ciphertext convolutions. Together, these innovations effectively mitigate the excessive computational burden of linear layers.

## F.2. Approximate Expression of Activation Layer

Scholars have found that, in order to maintain model inference accuracy over FHE, the use of bootstrapping operations must be carefully balanced. Too few bootstrapping operations can result in excessive noise accumulation, leading to significant inference error or decryption failure. Conversely, excessive bootstrapping significantly increases the overall inference time, making it another major bottleneck in FHE-based model inference. For FHE, the number of bootstrapping operations required is primarily determined by the multiplicative depth supported by the encryption parameters. Currently, the majority of multiplicative depth is consumed by non-linear activation functions, since the addition and multiplication alone cannot express non-linear transformations. A few researchers, such as Crawford et al. [38], have attempted to approximate activation functions using lookup tables for low-precision estimation. However, this often leads to considerable accuracy degradation, which prompts to adopt polynomial approximation technique instead. Gilad-Bachrach et al. [17] propose a square function to substitute the activation function, which yields satisfied results in a shallow five-layer network. However, they also point out several issues associated with square activation, such as gradient explosion during training, limited expressiveness, and hindered learning capacity. These problems would become more severe in deeper networks. Brutzkus et al. [39] also employ square activation in their LoLa system. Bourse et al. [40] propose an alternative approach by developing a method to compute nonlinear activation using sign functions. Lee et al. [21] empirically select a composite polynomial composed of minimax approximations of degree-5/15/27, which offers small approximation errors. However, this method consumes a large amount of multiplicative depth, requiring frequent bootstrapping operations

and leading to an inference time of approximately 3 hours under the environment: AMD Ryzen Threadripper PRO 3995WX at 2.096 GHz (64 cores) with 512 GB RAM. Park et al. [22] adopt a degree-2 Hermite polynomial to approximate the ReLU function. This design significantly reduces bootstrapping overhead due to its low multiplicative depth consumption. Nonetheless, the coarse-fined approximation results in a substantial drop in inference accuracy. Resorting to modeling the time overhead problem into a multi-objective optimization framework, Ao et al. [23] recently propose hierarchical mixed-degree polynomial activation to maximize accuracy while minimizing bootstrapping operations. However, their inference framework still faces a trade-off problem between accuracy and inference time, and the complex network architecture demands a high-performance deployment environment (at least 768 GB of RAM).

Towards this dilemma, our work proposes a novel activation function approximation strategy, aiming to optimize multiplicative depth consumption while improving the approximation quality of nonlinear activation functions. Specifically, a low-degree Legendre polynomial is used to approximate the smooth function SiLU, which has better representational capabilities. This approximation achieves a high level of inference accuracy while consuming only three levels of multiplicative depth. As a result, it significantly enhances the expressiveness of activation layers in the ciphertext domain, allowing the use of pretrained model parameters from plaintext environments. Our approach improves the overall model inference accuracy and effectively addresses the issue of excessive multiplicative depth consumption and poor performance on the approximation quality of existing low-degree polynomials.