

Bi-CryptoNets: Leveraging Different-Level Privacy For Encrypted Inference

Man-Jie Yuan, Zheng Zou, and Wei Gao^(✉)

National Key Laboratory for Novel Software Technology, Nanjing University, China
School of Artificial Intelligence, Nanjing University, China
{yuanmj, zouz, gaow}@lamda.nju.edu.cn

Abstract Privacy-preserving neural networks have attracted increasing attention in recent years, and various algorithms have been developed to keep the balance between accuracy, computational complexity and information security from the cryptographic view. This work takes a different view from the input data and structure of neural networks. We decompose the input data (e.g., some images) into sensitive and insensitive segments according to importance and privacy. The sensitive segment includes some important and private information such as human faces and we take strong homomorphic encryption to keep security, whereas the insensitive one contains some background and we add perturbations. We propose the bi-CryptoNets, i.e., plaintext and ciphertext branches, to deal with two segments, respectively, and ciphertext branch could utilize the information from plaintext branch by unidirectional connections. We adopt knowledge distillation for our bi-CryptoNets by transferring representations from a well-trained teacher neural network. Empirical studies show the effectiveness and decrease of inference latency for our bi-CryptoNets.

Keywords: Neural network · Encryption · Privacy-preserving inference.

1 Introduction

Recent years have witnessed increasing attention on privacy-preserving neural networks [14,4,25], which can be viewed as a promising security solution to the emerging Machine Learning as a Service (MLaaS) [30,13,15]. Specifically, some clients can upload their encrypted data to the powerful cloud infrastructures, and then obtain machine learning inference services; the cloud server performs inference without seeing clients' sensitive raw data by cryptographic primitives, and preserve data privacy.

Privacy-preserving neural networks are accompanied with heavy computational costs because of homomorphic encryption [12,2], and various algorithms and techniques have been developed to keep the balance between accuracy, information security and computation complexity. For example, Brutzkus et al. proposed the LoLa network for fast inference over a single image based on well-designed packing methods [4], and Lou and Jiang introduced the circuit-based network SHE with better accuracies [27], implemented with the TFHE scheme [7]. Dathathri et al. presented the compiler CHET to optimize data-flow graph of HE computations for privacy-preserving neural networks [10]. Yin et al. presented a comprehensive survey on privacy-preserving networks [39].

Previous studies mostly encrypted the entire input data and treated indiscriminately. In some applications, however, input data may consist of sensitive and insensitive segments according to different importance and privacy. As shown in Figure 1, a fighter is more sensitive than background such as sky and mountains in a military picture, and a human face is more private than landscape in a photo.



Figure 1: An illustration for input data (e.g., some images¹), which consists of two segments with different importance and privacy.

This work presents new privacy-preserving neural network from the view of input data and network structure, and the main contributions can be summarized as follows:

- We decompose input data (e.g., images) into sensitive and insensitive segments according to importance and privacy. We adopt strong homomorphic encryption to keep the security of sensitive segment, yet mingle some perturbations [11,38] to insensitive segment. This could reduce computational overhead and perform private inference of low latency, without unnecessary encryption on insensitive segment.
- We propose the bi-CryptoNets, i.e., ciphertext and plaintext branches, to deal with sensitive and insensitive segments, respectively. The ciphertext branch could use information from plaintext branch by unidirectional connections, but the converse direction does not hold because of the spread of ciphertexts. We integrate features for the final predictions, from the outputs of both branches.
- We present the feature-based knowledge distillation to improve the performance of our bi-CryptoNets, from a teacher of convolutional neural network trained on the entire data without decomposition.
- We present empirical studies to validate the effectiveness and decrease of inference latency for our bi-CryptoNets. We could improve inference accuracy by 0.2% ~ 2.1%, and reduce inference latency by $1.15\times \sim 3.43\times$ on a single image, and decrease the amortized latency by $4.59\times \sim 13.7\times$ on a batch of images.

The rest of this work is organized as follows: Section 2 introduces relevant work. Section 3 presents our bi-CryptoNets. Section 4 proposes the feature-based knowledge distillation. Section 5 conducts experiments, and Section 6 concludes with future work.

2 Relevant Work

Homomorphic encryption. Homomorphic Encryption (HE) is a cryptosystem that allows operations on encrypted data without requiring access to a secret key [12]. In addition to encryption function E and decryption function D , HE scheme provides two operators \oplus and \otimes such that, for every pair of plaintexts x_1 and x_2 ,

$$D(E(x_1) \oplus E(x_2)) = x_1 + x_2, \quad D(E(x_1) \otimes E(x_2)) = x_1 \times x_2,$$

where $+$ and \times are the standard addition and multiplication, respectively. Hence, we could directly perform private addition, multiplication and polynomial functions over encrypted data, by using \oplus and \otimes operators without knowing true values in plaintexts.

¹ Download from ILSVRC dataset [38].

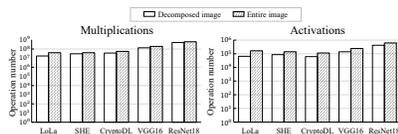


Figure 2: Counts of HE multiplications and activations for decomposed and entire image.

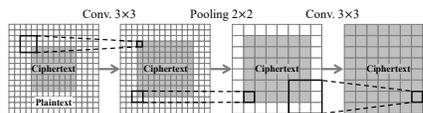


Figure 3: The fast spread of ciphertext via CNN layers.

The commonly-used HE cryptosystems include CKKS [6], BGV [3] and TFHE [7] for privacy-preserving machine learning. The BGV and TFHE schemes support integer and bits computations, while the CKKS scheme supports fixed-point computations. For CKKS, a ciphertext is an element in the ring \mathcal{R}_p^2 , where $\mathcal{R}_p = \mathbb{Z}_p[x]/(x^N + 1)$ is the residue ring of polynomial ring, with polynomial degree N . In addition to HE, a recent study introduced a new data encryption scheme by incorporating crucial ingredients of learning algorithm, specifically the Gini impurity in random forests [37].

Privacy-preserving neural networks. Much attention has been paid on the privacy-preserving neural networks in recent years. For example, Gilad-Bachrach et al. proposed the first CryptoNets to show the feasibility of private CNN inference by HE scheme [14]. Brutzkus et al. proposed the LoLa to optimize the implementation of matrix-vector multiplication by different packing methods, and achieve fast inference for single image [4]. In those studies, the ReLU activation has been replaced by squared activation because the used HE schemes merely support polynomial functions. Lou and Jiang presented a circuit-based quantized neural network SHE implemented by TFHE scheme, which could perform ReLU and max pooling by comparison circuits to obtain good accuracies [27]. Recent studies utilized the CKKS scheme to avoid quantization of network parameters for better accuracies [9,25].

Knowledge Distillation. Knowledge distillation focuses on transferring knowledge from a pretrained teacher model to a student model [22,26]. Traditional methods aimed to match the output distributions of two models by minimizing the Kullback-Leibler divergence loss [22]. Subsequently, a variety of innovative distillation techniques have been developed. Romero et al. proposed Fitnets [31], which matches the feature activations, while Zagoruyko and Komodakis suggested matching attention maps between two models [40]. Gou et al. gave an extensive review on knowledge distillation [17].

Threat model. Our threat model follows previous studies on private inference [1,28,8]. A honest but curious cloud-based machine learning service hosts a network, which is trained on plaintext data (such as public datasets) and hence the network weights are not encrypted [14]. To ensure the confidentiality of client’s data, the client could encrypt the sensitive segment of data by using HE scheme and send data to the cloud server for performing private inference service without decrypting data or accessing the client’s private key. Only the client could decrypt the inference result by using the private key.

3 Our bi-CryptoNets

In this section, we will present new privacy-preserving neural network according to different-level privacy of input data. Our motivation is to decompose input data into

sensitive and *insensitive* segments according to their privacy. The sensitive segment includes some important and private information such as human faces in an image, whereas the insensitive one contains some background information, which is not so private yet beneficial to learning algorithm.

Based on such decomposition, we could take the strongest homomorphic encryption to keep data security for sensitive segment, while mingle with some perturbations [11,38] to the data of insensitive segment. This is quite different from previous private inference [14,25], where homomorphic encryption is applied to the entire input data.

Notice that we can not directly use some previous neural networks to tackle such decomposition with much smaller computational cost, as shown in Figure 2. We compare with three private networks LoLa [4], SHE [27], CryptoDL [20], and two conventional networks VGG16 and ResNet18. Prior networks can not reduce HE operations because ciphertexts could quickly spread over the entire image via several convolution and pooling operations, as shown in Figure 3.

Our idea is quite simple and intuitive for the decomposition of input data. We construct a bi-branch neural network to deal with the sensitive and insensitive segments, which are called *ciphertext* and *plaintext* branch, respectively. The ciphertext branch can make use of features from plaintext branch by unidirectional connections, while the converse direction does not hold because of the quick spread of ciphertexts. We take the feature integration for final predictions from the outputs of two branches of network.

Figure 4 presents an overview for our new privacy-preserving neural network, and it is short for *bi-CryptoNets*. We will go to the details of *bi-CryptoNets* in the following.

3.1 The bi-branch of neural network

We construct ciphertext and plaintext branches to deal with the sensitive and insensitive segments of an input instance, respectively. The plaintext branch deals with the entire input instance, where the insensitive segment is mingled with some perturbations [11,38], while the sensitive segment is simply filled with zero. The ciphertext branch tackles the sensitive segment with the homomorphic encryption due to its privacy.

The HE operations are restricted in the ciphertext branch to avoid the ciphertext spread on plaintext branch. This could keep the proportion of HE operations in a stable level, rather than previous close approximation to 1 as depth increases [14]. Hence, our bi-branch structure could decrease HE operations. Moreover, plaintext branch can be computed with full precision for performing inference, rather than quantized homomorphic ciphertext in prior privacy-preserving neural networks [20].

3.2 The unidirectional connections

It is well-known that one great success of deep learning lies in the strong features or representations by plentiful co-adaptations of neural network [21,35,16]. Our bi-branch neural network reduces some co-adaptations of features obviously, to restrict the spread of ciphertexts. Hence, it is necessary to strengthen features' representations by exploiting some correlations in bi-branches.

Notice that the ciphertext branch can make use of features from plaintext branch, while the converse direction does not hold because of the spread of ciphertexts. In

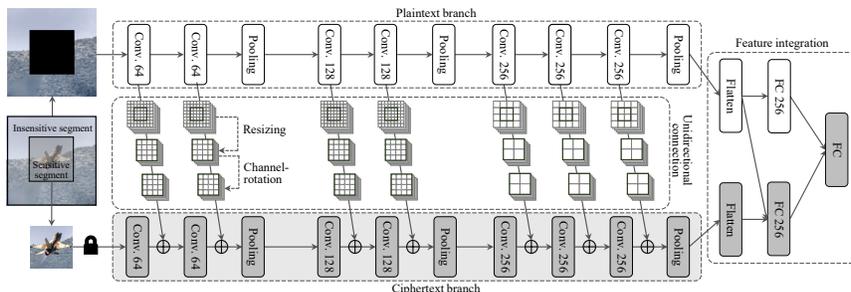


Figure 4: The overview of our proposed bi-CryptoNets, where grey layers are computed with encrypted inputs, and other layers are computed with plaintext inputs.

this work, we consider the simple addition to utilize features from plaintext branch as convolutional neural network [19,36,29]. Concatenation is another effective way to utilize features [24,23], while it yields more HE operations and computational costs.

One reason for addition is of small computational costs, since it yields relatively few homomorphic additions between plaintext and ciphertext, as shown in Table 1. We compare with addition (Add_{PC}) and multiplication (Mul_{PC}) between a plaintext and a ciphertext, as well as addition (Add_{CC}) and multiplication (Mul_{CC}) between two ciphertexts. We compare with four operations under three popular HE schemes: BGV [3], CKKS [6] and TFHE [7], and it is evident that the Add_{PC} has the smallest latency.

Table 1: The comparisons of latency (ms) for three schemes and four operations.

HE Scheme	Add_{PC}	Add_{CC}	Mul_{PC}	Mul_{CC}
BGV	0.049	0.077	2.055	3.379
CKKS	0.039	0.077	0.173	0.390
TFHE	56.03	256.8	1018	1585

For unidirectional connections, we first resize the feature map of plaintext branch to fit the addition with ciphertext branch, and then introduce a channel-rotation for ciphertext channel to extract some information from multiple plaintext channels. Specifically, suppose there are L layers in ciphertext and plaintext branches, and denote by $\mathcal{F}_{c,l}$ and $\mathcal{F}_{p,l}$ their respective l -th layer. For every $l \in [L]$, we consider the following two steps:

i) Resizing feature map of plaintext branch

We begin with the resizing function

$$\mathbf{y}_{p,l} = \text{Resize}_l(\mathbf{x}_{p,l}) : \mathbb{R}^{h_{p,l} \times w_{p,l} \times \text{ch}_{p,l}} \rightarrow \mathbb{R}^{h_{c,l} \times w_{c,l} \times \text{ch}_{p,l}}$$

where $\mathbf{x}_{p,l} \in \mathbb{R}^{h_{p,l} \times w_{p,l} \times \text{ch}_{p,l}}$ denotes the output feature maps from the l -th plaintext layer $\mathcal{F}_{p,l}$ with $h_{p,l}$ rows, $w_{p,l}$ columns and $\text{ch}_{p,l}$ channels, and $h_{c,l}$ and $w_{c,l}$ denote the number of row and column in ciphertext branch.

For resizing function, a simple yet effective choice is the cropping [32], which maintains the center features of size $h_{c,l} \times w_{c,l} \times \text{ch}_{p,l}$ because those features can capture more important information around the sensitive segment. We can also select pooling or convolution operations for resizing function [19].

ii) Channel-rotation for ciphertext channel

We now introduce the channel-rotation function as follows:

$$\mathbf{z}_{p,l} = \text{CRot}_l(\mathbf{y}_{p,l}) = \left[\text{CRot}_l^1(\mathbf{y}_{p,l}), \dots, \text{CRot}_l^i(\mathbf{y}_{p,l}), \dots, \text{CRot}_l^{\text{ch}_{c,l}}(\mathbf{y}_{p,l}) \right]$$

where $\text{ch}_{c,l}$ is number of channels in ciphertext branch, and $\text{CRot}_l^i(\mathbf{y}_{p,l})$ denotes the i -th channel of $\text{CRot}_l(\mathbf{x})$, that is,

$$\text{CRot}_l^i(\mathbf{y}_{p,l}) = \sum_{j=1}^{\text{ch}_{p,l}} \mathbf{W}_{\text{CRot},l}^{i,j} \cdot \mathbf{y}_{p,l}^j \quad \text{for } i \in [\text{ch}_{c,l}].$$

Here, $\mathbf{y}_{p,l}^j \in \mathbb{R}^{h_{c,l} \times w_{c,l}}$ denotes the j -th channel of $\mathbf{y}_{p,l}$ for $j \in [\text{ch}_{p,l}]$, and $\mathbf{W}_{\text{CRot},l} = (\mathbf{W}_{\text{CRot},l}^{i,j})_{\text{ch}_{c,l} \times \text{ch}_{p,l}}$ is the weight matrix for channel-rotation in the l -th layer.

The channel-rotation is helpful for ciphertext channels to extract useful information automatically from plaintext channels. This is because channel-rotation could compress, rotate and scale different channels of plaintext feature maps, without the requirements of channel-wise alignments and identical numbers and magnitudes of channels. Moreover, each channel, in connection with ciphertext branch, can be viewed as a linear combination of multiple plaintext channels, rather than one plaintext channel. Therefore, the ciphertext branch can get better information from plaintext branch and achieve better performance with the help of channel-rotation.

After resizing and channel-rotation, the ciphertext branch can make use of features from plaintext branch by addition. The unidirectional connections can be written as

$$\mathbf{x}_{c,l+1} = \mathcal{F}_{c,l+1}(\mathbf{x}_{c,l} + \mathbf{z}_{p,l}) = \mathcal{F}_{c,l+1}(\mathbf{x}_{c,l} + \text{CRot}_l(\text{Resize}_l(\mathbf{x}_{p,l}))),$$

where $\mathbf{x}_{c,l} \in \mathbb{R}^{h_{c,l} \times w_{c,l} \times \text{ch}_{c,l}}$ denotes the output feature maps of $\mathcal{F}_{c,l}$ with $h_{c,l}$ rows, $w_{c,l}$ columns and $\text{ch}_{c,l}$ channels.

Finally, it is feasible to improve computational efficiency by implementing plaintext and ciphertext branches in parallel, because computing the plaintext branch and unidirectional connections is much faster than that of ciphertext branch.

3.3 The feature integration

We now design a two-layer neural network to integrate the outputs from ciphertext and plaintext branches. In the first layer, we split all neurons into two halves, one half for ciphertext yet the other for plaintext. The plaintext neurons are only connected with the outputs from plaintext branch, whereas the ciphertext neurons have full connections. We take full connections in the second layer.

Specifically, there are n_1 (even) and n_2 neurons in the first and second layer, respectively. Let $\mathbf{x}_c \in \mathbb{R}^{n_c}$ and $\mathbf{x}_p \in \mathbb{R}^{n_p}$ denote the flattened outputs from ciphertext and plaintext branches with n_c and n_p features, respectively. Then, the outputs of ciphertext and plaintext neurons in the first layer can be given by, respectively,

$$\mathbf{x}_p^{(1)} = \sigma(\mathbf{W}'_{p,1} \mathbf{x}_p + \mathbf{b}'_1), \quad \mathbf{x}_c^{(1)} = \sigma(\mathbf{W}_{c,1} \mathbf{x}_c + \mathbf{W}_{p,1} \mathbf{x}_p + \mathbf{b}_1),$$

where $\sigma(\cdot)$ is an activation function, $\mathbf{b}_1, \mathbf{b}'_1 \in \mathbb{R}^{n_1/2}$ are bias vectors, $\mathbf{W}_{c,1} \in \mathbb{R}^{n_c \times n_1/2}$ and $\mathbf{W}_{p,1} \in \mathbb{R}^{n_p \times n_1/2}$ are the weight for ciphertext neurons, yet $\mathbf{W}'_{p,1} \in \mathbb{R}^{n_p \times n_1/2}$ is the weight for plaintext neurons. The final output in the second layer can be given by

$$\mathbf{x}_{\text{out}} = \sigma(\mathbf{W}_{c,2}\mathbf{x}_c^{(1)} + \mathbf{W}_{p,2}\mathbf{x}_p^{(1)} + \mathbf{b}_2),$$

where $\mathbf{b}_2 \in \mathbb{R}^{n_2}$ is a bias vector, and $\mathbf{W}_{c,2} \in \mathbb{R}^{n_1/2 \times n_2}$ and $\mathbf{W}_{p,2} \in \mathbb{R}^{n_1/2 \times n_2}$ are the weight matrices for ciphertext and plaintext neurons, respectively. Here, we consider two-layer neural network for simplicity, and similar constructions could be made for deeper neural networks based on the splitting of plaintext and ciphertext neurons.

4 Knowledge Distillation for bi-CryptoNets

Knowledge distillation has been an effective way to improve learning performance by distilling knowledge from a teacher network [22,17]. This section develops a feature-based knowledge distillation for bi-CryptoNets, as shown in Figure 5. The basic idea is to supplement the representations of two branches of bi-CryptoNets by imitating a teacher network.

For the teacher network, we learn a classical convolutional neural network from training data over the entire images without decomposition. Hence, the teacher network has plentiful connections between different neurons, and strengthens the intrinsic correlations for better performance.

We learn the representations of our two branches from the corresponding intermediate representations of teacher network, and also learn the final outputs of our bi-CryptoNets from that of teacher network. This is partially motivated from previous knowledge distillations on internal representations [31].

Specifically, we first pad the output of ciphertext branch with 0 to match the size of plaintext output and add them together, as in convolutional neural network literature [16,18]. We then train two branches of bi-CryptoNets to learn teacher’s intermediate representations by minimizing the following loss function:

$$\mathcal{L}_{\text{IR}}(\mathbf{W}_T^h, \mathbf{W}_c, \mathbf{W}_p) = \|\mathcal{F}_T^h(\mathbf{x}, \mathbf{W}_T^h) - (\text{Pad}(\mathcal{F}_c(\mathbf{x}, \mathbf{W}_c)) + \mathcal{F}_p(\mathbf{x}, \mathbf{W}_p))\|_2^2, \quad (1)$$

where $\mathcal{F}_T^h(\cdot; \mathbf{W}_T^h)$ is the first h layers of teacher network of parameter \mathbf{W}_T^h , and $\mathcal{F}_c(\cdot; \mathbf{W}_c)$ and $\mathcal{F}_p(\cdot; \mathbf{W}_p)$ denote the ciphertext and plaintext branches of parameter \mathbf{W}_c and \mathbf{W}_p , respectively, and $\text{Pad}(\cdot)$ is the zero-padding function.

Here, we denote by h the corresponding h -th layer in the teacher network that has the same output size as plaintext branch’s output. Such loss can help the sum of outputs from two branches to approximate the teacher’s intermediate output. In this manner, two branches could obtain stronger representations, and this makes it much easier for further learning in the following layers of our bi-CryptoNets.

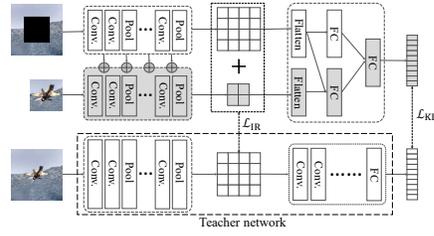


Figure 5: The overview of our feature-based knowledge distillation.

After training two branches, we then perform knowledge distillation to the whole network. The student network is trained such that its output is similar to that of the teacher and to the true labels. Given bi-CryptoNets’ output \mathbf{a}_S and teacher’s output \mathbf{a}_T , we first get softened output as

$$\mathbf{p}_T^\tau = \text{softmax}(\mathbf{a}_T/\tau) \quad \text{and} \quad \mathbf{p}_S^\tau = \text{softmax}(\mathbf{a}_S/\tau),$$

where $\tau > 1$ is a relaxation parameter. We try to optimize the following loss function:

$$\mathcal{L}_{\text{KD}}(\mathbf{W}_T, \mathbf{W}_{\text{BCN}}) = \mathcal{H}(\mathbf{y}_{\text{True}}, \mathbf{p}_s) + \lambda \mathcal{H}(\mathbf{p}_T^\tau, \mathbf{p}_s^\tau), \quad (2)$$

where \mathbf{y}_{True} is the true label, $\mathbf{p}_s = \text{softmax}(\mathbf{a}_S)$, \mathcal{H} refers to the cross-entropy, \mathbf{W}_T and \mathbf{W}_{BCN} are the parameters of teacher network and our bi-CryptoNets, respectively, and $\lambda \in [0, 1]$ is a hyperparameter balancing two cross-entropies. With our proposed method, two branches in our bi-CryptoNets can learn better representations, enhancing overall performance. More details can be found in Appendix G.

5 Experiments

This section presents empirical studies on datasets² MNIST, CIFAR-10 and CIFAR-100, which have been well-studied in previous private inference studies [4,27,25]. We develop different backbones for our bi-CryptoNets according to different datasets. We adopt the 3-layer CNN [14,4] as backbone for MNIST, and we take 11-layer CNN [5], VGG-16 [34] and ResNet-18 [19] for CIFAR-10 and CIFAR-100. Figure 4 presents our bi-CryptoNets with VGG-16 backbone, and more details are shown in Appendix C.

We take the CKKS scheme for sensitive segment with polynomial degree $N = 2^{14}$, and mingle with Gaussian noise for insensitive segment. We further improve packing method for our bi-CryptoNets to perform inference for multiple images simultaneously.

For simplicity, we focus on the regular images, where the area of sensitive segment is restricted to a quarter in the center of every image, and we could take some techniques of zooming, stretching and rotating to adjust those irregular images.

Ablation studies

We conduct ablation studies to verify the effectiveness of feature-based knowledge distillation and the structures of our bi-CryptoNets. Table 2 presents the details of experimental results of inference accuracies on three datasets.

We first exploit the influence of feature-based knowledge distillation in our bi-CryptoNets. We consider two different methods: bi-CryptoNets without knowledge distillation and bi-CryptoNets with conventional knowledge distillation. As can be seen from Table 2, it is observable that our feature-based knowledge distillation could effectively improve the inference accuracy by 0.54% \sim 8.39%, in comparison to bi-CryptoNets without knowledge distillation; it also achieves better inference accuracy by 1.45% \sim 10.27% than that of conventional knowledge distillation, which fails to stably improve accuracy with our bi-branch structure.

² Downloaded from yann.lecun.com/exdb/mnist and www.cs.toronto.edu/~kriz/cifar

Table 2: The accuracies (%) on MNIST, CIFAR-10 and CIFAR-100 datasets (KD refers to the conventional knowledge distillation; FKD refers to our feature-based knowledge distillation).

Models	Knowledge	MNIST	CIFAR-10			CIFAR-100	
	Distillation	CNN-3	CNN-11	VGG-16	ResNet-18	VGG-16	ResNet-18
Backbone network (w/o decomposition)	w/o KD	99.21	90.99	93.42	94.30	72.23	74.00
bi-CryptoNets	w/o KD	98.60	81.91	90.36	92.04	64.85	69.46
bi-CryptoNets	KD	98.61	80.03	88.91	92.43	64.96	65.67
bi-CryptoNets (w/o uni. connections)	FKD	98.89	84.69	91.78	91.73	70.61	70.29
bi-CryptoNets (w/o channel-rotations)	FKD	98.90	90.09	92.14	91.86	71.43	71.40
bi-CryptoNets	FKD	99.15	90.30	93.27	93.91	72.35	73.33

We then study the influence of unidirectional connections in bi-CryptoNets, and consider two variants: our bi-CryptoNets without unidirectional connections and bi-CryptoNets with resizing yet without channel-rotations. As can be seen from Table 2, the unidirectional connections are helpful for ciphertext branch to extract useful information from plaintext branch. This is because the unidirectional connections with only resizing can enhance inference accuracy by 0.13% \sim 5.40%, and it could further improve accuracy by 0.26% \sim 5.61% with both resizing and channel-rotations.

We finally compare the inference accuracies of our bi-CryptoNets with the backbone network, which is trained and tested on the entire input images without decomposition. From Table 2, it is observable that our bi-CryptoNets with feature-based knowledge distillation achieves comparable or even better inference accuracies than the corresponding backbone networks without data decomposition; therefore, our proposed bi-CryptoNets could effectively compensate for the information loss of plaintext and ciphertext branches under the help of feature-based knowledge distillation.

Experimental comparisons

We compare our bi-CryptoNets with the state-of-the-art schemes on privacy-preserving neural networks [8,20,4,27,9,25]. We also implement the structure of backbone network with square activation for fair comparisons. We take the batch size 20 and 32 for MNIST and CIFAR-10, respectively. We employ commonly-used criteria in experiments, i.e., inference accuracy, inference latency and the number of homomorphic operations as in [8,28]. We also consider the important amortized inference latency in a batch of images [25], and concern the number of activations for ciphertexts [27,4,9]. Tables 3 and 4 summarize experimental comparisons on MNIST and CIFAR-10, respectively. Similar results could be made for CIFAR-100, presented in Appendix E.

From Table 3, our bi-CryptoNets can reduce HE operations by $2.35\times$ and inference latency by $3.43\times$ in contrast to backbone network. Our bi-CryptoNets could decrease the amortized latency by $13.7\times$, and achieve better accuracy (about 0.2%) than backbone because of our knowledge distillation and precise inference in the plaintext branch.

From Table 4, our bi-CryptoNets reduces HE operations by $1.43\times$, and inference latency by $1.15\times$ in contrast to backbone network. Our bi-CryptoNets could decrease the amortized latency by $4.59\times$ and improve accuracy by 2.1%. We also implement our

Table 3: The experimental comparisons on MNIST.

Scheme	HEOPs	Add _{CC}	Mul _{PC}	Act _C	Latency(s)	Amortized Latency(s)	Acc(%)
FCryptoNets [8]	63K	38K	24K	945	39.1	2.0	98.71
CryptoDL [20]	4.7M	2.3M	2.3M	1.6K	320	16.0	99.52
LoLa [4]	573	393	178	2	2.2	2.2	98.95
SHE [27]	23K	19K	945	3K	9.3	9.3	99.54
EVA [9]	8K	4K	4K	3	121.5	121.5	99.05
VDSCNN [25]	4K	2K	2K	48	105	105	99.19
Backbone CNN-3	1962	973	984	5	7.2	1.4	98.95
bi-CryptoNets (CNN-3)	830	406	418	6	2.1	0.1	99.15

Table 4: The experimental comparisons on CIFAR-10.

Scheme	HEOPs	Add _{CC}	Mul _{PC}	Act _C	Latency(s)	Amortized Latency(s)	Acc(%)
FCryptoNets [8]	701M	350M	350M	64K	22372	1398	76.72
CryptoDL [20]	2.4G	1.2G	1.2G	212K	11686	731	91.50
LoLa [4]	70K	61K	9K	2	730	730	76.50
SHE [27]	4.4M	4.4M	13K	16K	2258	2258	92.54
EVA [9]	135K	67K	67K	9	3062	3062	81.50
VDSCNN [25]	18K	8K	9K	752	2271	2271	91.31
Backbone CNN-11	1.0M	493K	521K	246	1823	228	88.21
bi-CryptoNets (CNN-11)	709K	341K	368K	246	1587	49	90.30
bi-CryptoNets (VGG-16)	3.4M	1.5M	1.9M	1.1K	2962	92	93.27
bi-CryptoNets (ResNet-18)	15.5M	6.9M	8.6M	2.8K	6760	211	93.91

bi-CryptoNets with VGG-16 and ResNet-18 as backbones, and deeper networks yield better inference accuracies, i.e., 93.27% for VGG-16 and 93.91% for ResNet-18.

From Tables 3 and 4, our bi-CryptoNets takes a good balance between inference accuracies and computation cost in comparison with other neural networks. Our bi-CryptoNets achieves lower inference latency and amortized latency than other methods except for LoLa, where light-weight neural network is implemented with complicated packing method and smaller HE operations. Our bi-CryptoNets takes relatively-good inference accuracies except for SHE, CryptoDL and VDSCNN, where deeper neural networks are adopted with larger computational overhead. Based on deeper backbones such as VGG-16 and ResNet-18, our bi-CryptoNets could achieve better inference accuracy, comparable inference latency and smaller amortized latency.

6 Conclusion

Numerous privacy-preserving neural networks have been developed to keep the balance between accuracy, efficiency and security. We take a different view from the input data and network structure. We decompose the input data into sensitive and insensitive segments, and propose the bi-CryptoNets, i.e., plaintext and ciphertext branches, to deal with two segments, respectively. We also introduce feature-based knowledge

distillation to strengthen the representations of our network. Empirical studies verify the effectiveness of our bi-CryptoNets. An interesting future work is to exploit multiple levels of privacy of input data and multiple branches of neural network, and it is also interesting to generalize our idea to other settings such as multi-party computation.

Acknowledgements

The authors want to thank the reviewers for their helpful comments and suggestions. This research was supported by National Key R&D Program of China (2021ZD0112802), NSFC (62376119) and CAAI-Huawei MindSpore Open Fund.

References

1. Boemer, F., Costache, A., Cammarota, R., Wierzynski, C.: nGraph-HE2: A high-throughput framework for neural network inference on encrypted data. In: WAHC@CCS. pp. 45–56 (2019)
2. Boulemtafes, A., Derhab, A., Challal, Y.: A review of privacy-preserving techniques for deep learning. *Neurocomputing* **384**, 21–45 (2020)
3. Brakerski, Z., Gentry, C., Vaikuntanathan, V.: (leveled) fully homomorphic encryption without bootstrapping. *ACM Trans. Comput. Theory* **6**(3), 1–36 (2014)
4. Brutzkus, A., Gilad-Bachrach, R., Elisha, O.: Low latency privacy preserving inference. In: ICML. pp. 812–821 (2019)
5. Chabanne, H., Wargny, A., Milgram, J., Morel, C., Prouff, E.: Privacy-preserving classification on deep neural network. *Cryptol. ePrint Arch.* (2017)
6. Cheon, J., Kim, A., Kim, M., Song, Y.: Homomorphic encryption for arithmetic of approximate numbers. In: ASIACRYPT. pp. 409–437 (2017)
7. Chillotti, I., Gama, N., Georgieva, M., Izabachène, M.: TFHE: Fast fully homomorphic encryption over the torus. *J. Cryptol.* **33**(1), 34–91 (2020)
8. Chou, E., Beal, J., Levy, D., Yeung, S., Haque, A., Fei-Fei, L.: Faster CryptoNets: Leveraging sparsity for real-world encrypted inference. *CoRR/abstract* **1811.09953** (2018)
9. Dathathri, R., Kostova, B., Saarikivi, O., Dai, W., Laine, K., Musuvathi, M.: EVA: An encrypted vector arithmetic language and compiler for efficient homomorphic computation. In: PLDI. pp. 546–561 (2020)
10. Dathathri, R., Saarikivi, O., Chen, H., Laine, K., Lauter, K., Maleki, S., Musuvathi, M., Mytkowicz, T.: CHET: An optimizing compiler for fully-homomorphic neural-network inferencing. In: PLDI. pp. 142–156 (2019)
11. Dwork, C., McSherry, F., Nissim, K., Smith, A.: Calibrating noise to sensitivity in private data analysis. *J. Priv. Confidentiality* **7**(3), 17–51 (2016)
12. Gentry, C.: Fully homomorphic encryption using ideal lattices. In: STOC. pp. 169–178 (2009)
13. Ghodsi, Z., Jha, N., Reagen, B., Garg, S.: Circa: Stochastic ReLUs for private deep learning. In: NeurIPS. pp. 2241–2252 (2021)
14. Gilad-Bachrach, R., Dowlin, N., Laine, K., Lauter, K., Naehrig, M., Wernsing, J.: Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In: ICML. pp. 201–210 (2016)
15. Gong, X., Chen, Y., Yang, W., Mei, G., Wang, Q.: InverseNet: Augmenting model extraction attacks with training data inversion. In: IJCAI. pp. 2439–2447 (2021)

16. Goodfellow, I., Bengio, Y., Courville, A.: *Deep Learning*. MIT Press (2016)
17. Gou, J., Yu, B., Maybank, S., Tao, D.: Knowledge distillation: A survey. *Int. J. Comput. Vis.* **129**(6), 1789–1819 (2021)
18. Hashemi, M.: Enlarging smaller images before inputting into convolutional neural network: Zero-padding vs. interpolation. *J. Big Data* **6**(1), 1–13 (2019)
19. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *CVPR*. pp. 770–778 (2016)
20. Hesamifard, E., Takabi, H., Ghasemi, M.: Deep neural networks classification over encrypted data. In: *CODASPY*. pp. 97–108 (2019)
21. Hinton, G., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Improving neural networks by preventing co-adaptation of feature detectors. *CoRR/abstract* **1207.0580** (2012)
22. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. *CoRR/abstract* **1503.02531** (2015)
23. Huang, G., Liu, Z., Maaten, L., Weinberger, K.: Densely connected convolutional networks. In: *CVPR*. pp. 2261–2269 (2017)
24. Iandola, F., Moskewicz, M., Ashraf, K., Han, S., Dally, W., Keutzer, K.: SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size. *CoRR/abstract* **1602.07360** (2016)
25. Lee, E., Lee, J., Lee, J., Kim, Y., Kim, Y., No, J., Choi, W.: Low-complexity deep convolutional neural networks on fully homomorphic encryption using multiplexed parallel convolutions. In: *ICML*. pp. 12403–12422 (2022)
26. Li, Z., Li, X., Yang, L., Zhao, B., Song, R., Luo, L., Li, J., Yang, J.: Curriculum temperature for knowledge distillation. In: *AAAI*. pp. 1504–1512 (2023)
27. Lou, Q., Jiang, L.: SHE: A fast and accurate deep neural network for encrypted data. In: *NeurIPS*. pp. 10035–10043 (2019)
28. Lou, Q., Lu, W., Hong, C., Jiang, L.: Falcon: Fast spectral inference on encrypted data. In: *NeurIPS*. pp. 2364–2374 (2020)
29. Radosavovic, I., Kosaraju, R., Girshick, R., He, K., Dollár, P.: Designing network design spaces. In: *CVPR*. pp. 10425–10433 (2020)
30. Ribeiro, M., Grolinger, K., Capretz, M.: MLaaS: Machine learning as a service. In: *ICMLA*. pp. 896–902 (2015)
31. Romero, A., Ballas, N., Kahou, S., Chassang, A., Gatta, C., Bengio, Y.: Fitnets: Hints for thin deep nets. In: *ICLR* (2015)
32. Ronneberger, O., Fischer, P., Brox, T.: U-Net: Convolutional networks for biomedical image segmentation. In: *MICCAI*. pp. 234–241 (2015)
33. Microsoft SEAL (release 4.0). <https://github.com/Microsoft/SEAL> (Mar 2022), microsoft Research, Redmond, WA.
34. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. *CoRR/abstract* **1409.1556** (2014)
35. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: A simple way to prevent neural networks from overfitting. *JMLR* **15**(1), 1929–1958 (2014)
36. Tan, M., Le, Q.: EfficientNet: Rethinking model scaling for convolutional neural networks. In: *ICML*. pp. 6105–6114 (2019)
37. Xie, X.R., Yuan, M.J., Bai, X.T., Gao, W., Zhou, Z.H.: On the Gini-impurity preservation for privacy random forests. In: *NeurIPS* (2023)
38. Yang, K., Yau, J., Fei-Fei, L., Deng, J., Russakovsky, O.: A study of face obfuscation in imagenet. In: *ICML*. pp. 25313–25330 (2022)
39. Yin, X., Zhu, Y., Hu, J.: A comprehensive survey of privacy-preserving federated learning: A taxonomy, review, and future directions. *ACM Comput. Surv.* **54**(6), 131:1–131:36 (2021)
40. Zagoruyko, S., Komodakis, N.: Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In: *ICLR* (2017)

Supplementary Material (Appendix)

A CKKS Scheme

In this section, we will describe CKKS scheme in more details. CKKS is a HE scheme that supports fixed-point arithmetic operations on encrypted data. The ciphertext of the CKKS scheme is an element $\mathbf{c} \in \mathcal{R}^2$, where \mathcal{R} denotes the polynomial ring $\mathcal{R}_p^N = \mathbb{Z}_p[x]/(x^N + 1)$, $\mathbb{Z}_p[x] = \mathbb{Z}/p\mathbb{Z}$, and $p = \prod_{i=1}^l q_i$ is a product of l prime numbers, N is the polynomial degree. In CKKS scheme, we can encrypt $n = N/2$ plaintext in $N/2$ slots of a single ciphertext, and perform Single Instruction Multiple Data (SIMD) operations at no extra cost. Let $E(\cdot)$ and $D(\cdot)$ denote the encryption and decryption function respectively. Let $\mathbf{c}, \mathbf{c}_1, \mathbf{c}_2, \mathbf{c}_{\text{add}}, \mathbf{c}_{\text{mult}}, \mathbf{c}_{\text{rot}}$ denote ciphertext, and \mathbf{u} denotes a plaintext vector in \mathbb{R}^n . CKKS scheme supports following homomorphic operations:

- Homomorphic addition:

$$\text{Add}(\mathbf{c}, \mathbf{u}) \rightarrow \mathbf{c}_{\text{add}}, \text{ s.t. } D(\mathbf{c}_{\text{add}}) = D(\mathbf{c}_1) + \mathbf{u},$$

$$\text{Add}(\mathbf{c}_1, \mathbf{c}_2) \rightarrow \mathbf{c}_{\text{add}}, \text{ s.t. } D(\mathbf{c}_{\text{add}}) = D(\mathbf{c}_1) + D(\mathbf{c}_2).$$

- Homomorphic multiplication:

$$\text{Mult}(\mathbf{c}, \mathbf{u}) \rightarrow \mathbf{c}_{\text{mult}}, \text{ s.t. } D(\mathbf{c}_{\text{mult}}) = D(\mathbf{c}_1) \cdot \mathbf{u},$$

$$\text{Mult}(\mathbf{c}_1, \mathbf{c}_2) \rightarrow \mathbf{c}_{\text{mult}}, \text{ s.t. } D(\mathbf{c}_{\text{mult}}) = D(\mathbf{c}_1) \cdot D(\mathbf{c}_2),$$

where \cdot is the element-wise multiplication between two vectors.

- Homomorphic rotation:

$$\text{Rot}(\mathbf{c}, r) \rightarrow \mathbf{c}_{\text{rot}}, \text{ s.t. } D(\mathbf{c}_{\text{rot}}) = \langle D(\mathbf{c}) \rangle_r,$$

where $r \in \mathbb{Z}$, $\langle \cdot \rangle_r$ denotes cyclically shifting a vector by r to the left when $r > 0$, and shifting by $-r$ to the right when $r < 0$.

Gilad-Bachrach et al. performed single instruction multiple data operations without additional costs by packing multiple plaintext into one ciphertext for HE [14]. Specifically, a vector of $N/2$ fixed-point numbers can be packed into $N/2$ slots of a single ciphertext for CKKS [6], and element-wise multiplication and addition between two vectors can be performed with only one HE operation. Such technique has been used to reduce HE computations in privacy-preserving neural networks, such as packing multiple channels of input data or multiple neurons into one ciphertext [4,9].

B Experiment Details

Datasets. We conduct experiments on MNIST, CIFAR-10 and CIFAR-100 datasets. MNIST has 10 output classes with 60,000 training images and 10,000 testing images, and each image’s size is $28 \times 28 \times 1$. CIFAR-10 and CIFAR-100 have 10 and 100 output classes respectively, and both of them consist of 50,000 training images and 10,000 testing images, and each image’s size is $32 \times 32 \times 3$.

Parameter settings. Datasets have been preprocessed by image centering, that is, subtracting mean and dividing standard deviation, and images have been augmented by using the random horizontal flips and random crops. We also adopt the squared activation in ciphertext branch and feature integration, since CKKS only supports homomorphic addition and multiplication. In training, we use Adam optimizer to train our network with batch size of 256. We adopt the learning rate of 0.0001 with cosine decay. For feature-based knowledge distillation, the teacher network takes backbone network with ReLU activation and hyper-parameters $\tau = 4$ and $\lambda = 0.9$. Our experiments are conducted on the HE scheme library SEAL [33] on AMD Ryzen Threadripper 3970X at 2.2 GHz with 256GB of RAM, running the Ubuntu 20.04 operating system.

State-of-the-art privacy-preserving neural networks.

- FCryptoNets [8]: A privacy-preserving neural network that leverages sparse representations in the underlying cryptosystem to accelerate inference. It develops a pruning and quantization approach that achieves maximally-sparse encodings and minimizes approximation error. Such approach can reduce the number of HE operations, and achieve short inference latency.
- CryptoDL [20]: A privacy-preserving neural network with polynomial activations. It use low degree polynomials to approximate ReLU activations, and implement deeper networks, achieving better performance in accuracy.
- LoLa [4]: A privacy-preserving neural network with well-designed packing methods. It can pack one image into several ciphertext, and achieve fast inference for one single image with much fewer HE operations. On the other hand, it can only perform inference for one image every time.
- SHE [27]: A privacy-preserving neural network implemented with TFHE scheme. It use boolean circuit to implement neural network by quantizing network weights into powers of 2, and it can support ReLU activation and max pooling. Therefore, it can reach higher accuracy.
- EVA [9]: A privacy-preserving neural network that implemented by a compiler for HE with CKKS scheme. The compiler can optimize the hyperparameters and the computation graph for HE. Therefore, it can decrease the HE computation depth and shorten the inference latency for private-preserving neural networks.
- VDSCNN [25]: A privacy-preserving neural network implemented with CKKS scheme. It uses polynomial activations and it also packs multiple channels of the image into one ciphertext. Therefore, it can perform inference with fewer HE additions and multiplications.

C Network Architectures

In this section, we will present the structure for other variants of our bi-CryptoNets in details.

Figure 6 presents the overview structure of our proposed bi-CryptoNets that takes 11-layer CNN as backbone. Each branch has the following layout: (i) 2 layers of 3×3 convolution with stride of (1, 1), padding of (1, 1) and 32 output maps; (ii) 2×2 average

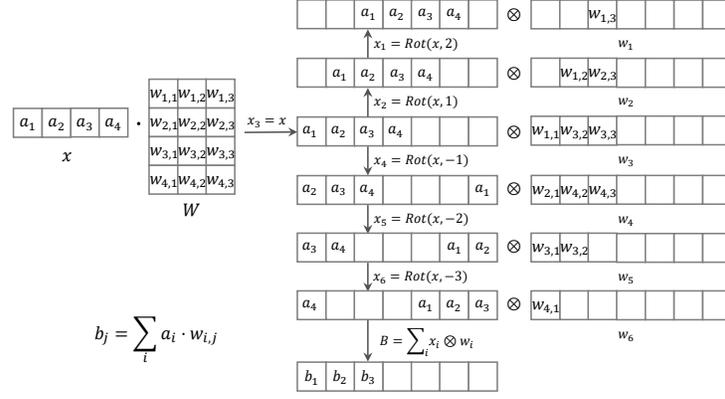


Figure 8: Homomorphic vector-matrix multiplication.

and 128 output maps; (vi) skip connection by a 3×3 convolution with stride of (2, 2) and 128 output maps with (iii) as input, and adding the results with (v); (vii) 3×3 convolution with stride of (2, 2) and 256 output maps; (viii) 3×3 convolution with stride of (1, 1), padding of (1, 1) and 256 output maps; (ix) skip connection by a 3×3 convolution with stride of (2, 2) and 256 output maps with (vi) as input, and adding the results with (viii); (x) 3×3 convolution with stride of (2, 2) and 512 output maps; (xi) 3×3 convolution with stride of (1, 1), padding of (1, 1) and 512 output maps; (xii) skip connection by a 3×3 convolution with stride of (2, 2) and 512 output maps with (viii) as input, and adding the results with (xi). We build unidirectional connections after each convolutional layers in the plaintext branch. After two branches, we have a 2-layer feature integration: In the first layer, we have two fully connected layers with 256 plaintext outputs and 256 ciphertext outputs, respectively; in the second layer, we have fully connected layer with 10 ciphertext outputs.

D Implementation Details

In this section, we will introduce our implementation for bi-CryptoNets with CKKS scheme in more details. A privacy-preserving neural network mainly consists of three types of layers: fully connected layer, convolutional layer and pooling layer.

Homomorphic fully connected layer

For an input $\mathbf{x} \in \mathbb{R}^{n_1}$, weight matrix $W \in \mathbb{R}^{n_1 \times n_2}$ and bias vector $\mathbf{k} \in \mathbb{R}^{n_2}$, the fully connected layer can be written as

$$\mathbf{y} = \sigma(\mathbf{x}^T W + \mathbf{k}),$$

where $\sigma(\cdot)$ is the activation function.

After adding them together, it will finally be multiplied by a mask vector, to only keep the values in the slots we need. We also can compute multiple convolution kernels for multiple input channels in parallel.

Homomorphic pooling layer

Similar to prior works in private inference, we adopt sum pooling in our network. The implementation of pooling layer is similar to that of homomorphic convolutional layer. Nevertheless, we do not need to multiply weights in convolution kernels.

BHW Packing

For inference, we could implement our bi-CryptoNets by packing method for speed on multiple images simultaneously, motivated by [10]. We introduce the new BHW packing for our bi-CryptoNets without extra costs, and the basic idea is to pack multiple segments into one ciphertext for decomposed images. Specifically, we pack n of $h \times w$ matrices into one ciphertext for a batch of n images and their sensitive segments with channel c , height h and width w . It requires c ciphertext for n images with our BHW packing.

Batch packing [14,8] presents an efficient method for a large number of images, to pack the same pixels in a batch of images into one ciphertext. Another relevant HW packing tries to pack one channel of pixels in an image into one ciphertext [10,25], which could greatly reduce HE computations and inference latency for single image.

Figure 10 presents a simple demonstration for three packing methods. Our

BHW packing could achieve low inference latency as that of HW packing over one single image, and take comparable inference latency as that of batch packing on multiple images of smaller batch size. This is because our BHW packing method could pack more segments into one ciphertext for the decomposed images, and only small sensitive segments are required to be encrypted.

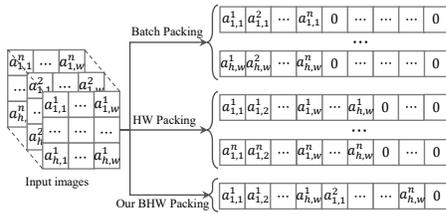


Figure 10: A simple demonstration for three packing methods.

E Results on CIFAR-100

In this section, we will present our experiment results on CIFAR-100 dataset. As shown in Table 5, our bi-CryptoNets can achieve similar results on CIFAR-100 as on CIFAR-10 dataset, since the network structure are similar to each other. Our bi-CryptoNets with CNN as backbone can also decrease the amortized latency by $4.67\times$ and improve accuracy by 2.48% in comparison with backbone network. Our bi-CryptoNets with VGG-16 and ResNet-18 as backbones still can achieve better inference accuracy, comparable inference latency and smaller amortized latency.

Table 5: The experimental comparisons on CIFAR-100, and abbreviations are similar to Table 3.

Scheme	HEOPs	Add _{CC}	Mul _{PC}	Act _C	Latency(s)	Amortized Latency(s)	Acc(%)
VDSCNN [25]	18K	8K	9K	752	3942	3942	69.43
Backbone CNN-11	1.0M	493K	521K	246	1829	229	66.30
bi-CryptoNets with CNN-11	709K	341K	368K	246	1588	49	68.78
bi-CryptoNets with VGG-16	3.4M	1.5M	1.9M	1.1K	2969	93	72.35
bi-CryptoNets with ResNet-18	15.5M	6.9M	8.6M	2.8K	6765	212	73.33

Table 6: The accuracies(%) on MNIST, CIFAR-10 and CIFAR-100 datasets.

Models	MNIST	CIFAR-10			CIFAR-100	
	CNN-3	CNN-11	VGG-16	ResNet-18	VGG-16	ResNet-18
Backbone network without decomposition	99.21	90.99	93.42	94.30	72.35	74.00
bi-CryptoNets (Convolution resizing)	98.96	89.69	92.24	92.05	71.51	70.79
bi-CryptoNets (Cropping resizing)	99.15	90.30	93.27	93.91	72.23	73.33

F Unidirectional Connections

In Section 3.2, we use cropping as resizing function for simplicity. Another way is to use a convolutional layer with strides as resizing function. In this section, we will show that cropping is a more effective way.

Similar to the skip connections of ResNet, we can simply use a convolutional layer as resizing functions. For example, if both the height and width of the sensitive segment is a half of that of insensitive segment, we can use a 3×3 convolution with stride of $(2, 2)$ as resizing functions. However, such approach has following disadvantages.

First, extra convolutional layers would introduce more parameters, which would make the network harder to train, and it would also lead to larger computation overhead. Second, when using cropping as resizing functions, it can bring representation from the plaintext branch at similar depth to the corresponding layers in ciphertext branch, whereas convolutions will bring higher orders representation to the layers in ciphertext branch. Intuitively, it is more natural to adding representations from similar depth together, compared with adding a high-order representation with low-order representation. Finally, cropping is an easier and more flexible way for resizing feature maps with various shapes, compared with convolution. Therefore, we simply use cropping as resizing functions in this paper.

Moreover, to compare the effectiveness of two resizing functions, we conduct experiments on MNIST, CIFAR-10 and CIFAR-100 with different network structures. We replace the cropping resizing function in our bi-CryptoNets with 3×3 convolution by stride of $(2, 2)$. We train both networks with feature-based knowledge distillation, and compare their accuracy. As shown in table 6, using cropping as resizing function can

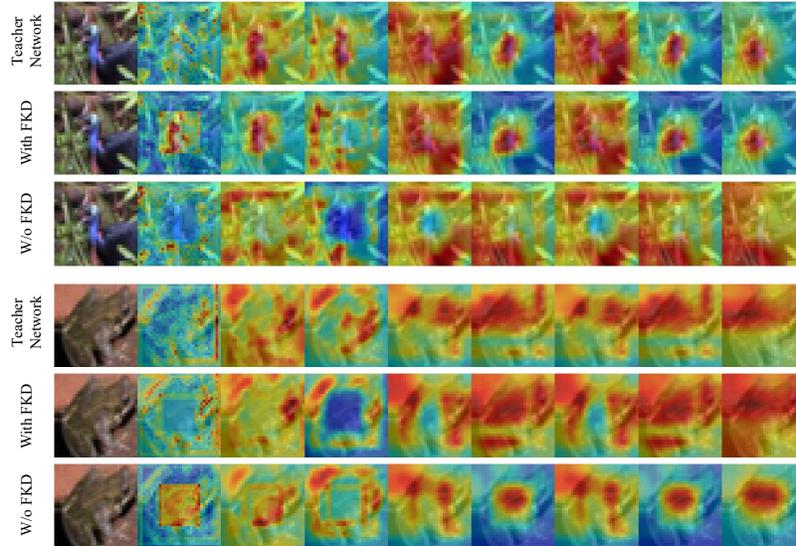


Figure 11: Attention maps in different layers of teacher network (CNN-11), bi-CryptoNets with feature-based knowledge distillation and bi-CryptoNets without knowledge distillation on CIFAR-10 images.

reach better accuracies in all datasets and with all variants of bi-CryptoNets structures, which is consistent with our analysis.

G Visualization for the effectiveness of feature-based knowledge distillation

In this section, we will show the effectiveness of feature-based knowledge distillation for bi-CryptoNets by drawing the attention maps at each layer. We compare our bi-CryptoNets trained by feature-based knowledge distillation with the teacher network and the bi-CryptoNets trained without knowledge distillation. Figure 11 presents the attention maps of each layer for bi-CryptoNets that takes 11-layer CNN as backbone on CIFAR-10 images, and Figure 11 presents the attention maps of each layer for bi-CryptoNets that takes VGG-16 as backbone. We can observe that the attention maps of bi-CryptoNets with feature-based knowledge distillation is very close to those of teachers, whereas bi-CryptoNets without knowledge distillation fails to capture the important patterns in the images, resulting the wrong predictions. Therefore, our feature-based knowledge distillation is an effective way to improve the performance of bi-CryptoNets, so that our bi-CryptoNets can achieve comparable performance as teacher network.

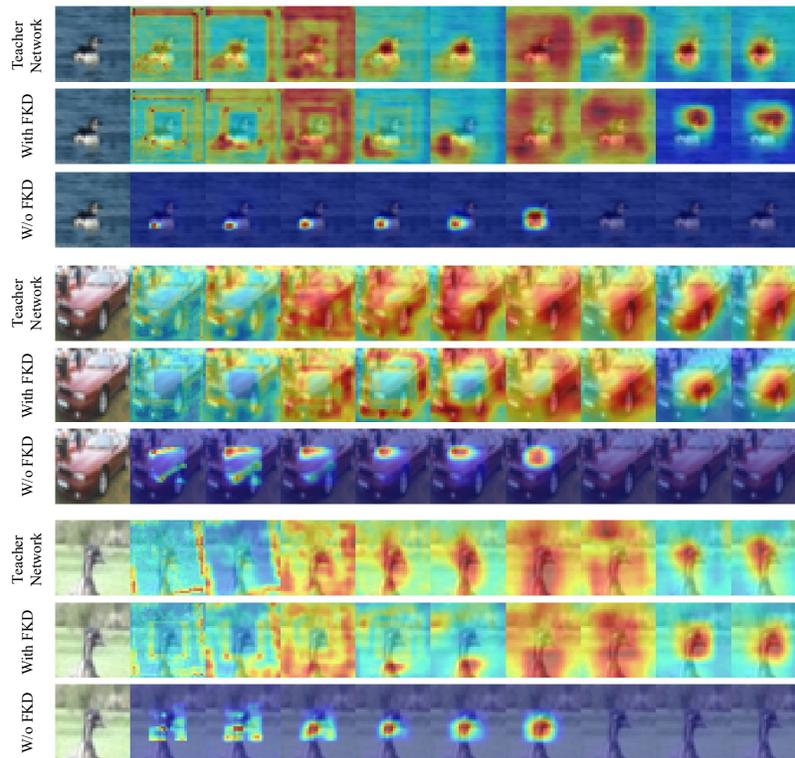


Figure 12: Attention maps in different layers of teacher network (VGG-16), bi-CryptoNets with feature-based knowledge distillation and bi-CryptoNets without knowledge distillation on CIFAR-10 images.