

A Comparative Study of Image Disguising Methods for Confidential Outsourced Learning

Sagar Sharma

Bytedance

Seattle, WA

sagar.sharma@bytedance.com

Yuechun Gu, Keke Chen

Trustworthy and Intelligent Computing Lab

Marquette University, Milwaukee WI

{ethan.gu, keke.chen}@marquette.edu

Abstract

Large training data and expensive model tweaking are standard features of deep learning for images. As a result, data owners often utilize cloud resources to develop large-scale complex models, which raises privacy concerns. Existing solutions are either too expensive to be practical or do not sufficiently protect the confidentiality of data and models. In this paper, we study and compare novel *image disguising* mechanisms, DisguisedNets and InstaHide, aiming to achieve a better trade-off among the level of protection for outsourced DNN model training, the expenses, and the utility of data. DisguisedNets are novel combinations of image blocktization, block-level random permutation, and two block-level secure transformations: random multidimensional projection (RMT) and AES pixel-level encryption (AES). InstaHide is an image mixup and random pixel flipping technique [16]. We have analyzed and evaluated them under a multi-level threat model. RMT provides a better security guarantee than InstaHide, under the Level-1 adversarial knowledge with well-preserved model quality. In contrast, AES provides a security guarantee under the Level-2 adversarial knowledge, but it may affect model quality more. The unique features of image disguising also help us to protect models from model-targeted attacks. We have done an extensive experimental evaluation to understand how these methods work in different settings for different datasets.

I. INTRODUCTION

Deep Neural Networks (DNN) have shown impressive performance across diverse domains such as image classification, natural language processing, speech recognition, and recommendation systems. However, DNN training is resource-intensive and time-consuming, requiring large training data, careful model architecture selection, and exhaustive model parameter tweaking. As a result, data owners or model developers often utilize multiple cloud GPUs or online model training services, such as Google Colab, to lower their costs.

Despite its popularity, outsourcing DNN learning to the cloud raises privacy and security concerns about the sensitive training data and trained models [31], [5]. On the one hand, cloud users cannot verifiably prevent the cloud provider from getting access to their data. In practice, using public clouds often means fully trusting your cloud provider. On the other hand, public cloud providers are not immune to

security attacks, which may lead to data breaches through insider [4], [5] and external attacks [25], [37]. Additionally, membership inference attacks [34], model inversion attacks [10], and adversarial example exploration [3], [28] can be applied to models directly to explore the training examples in DNN learning. Therefore, data and models in training, testing, and transferring between the cloud and the client are seriously threatened.

There have been a few efforts trying to address this critical issue. However, all of them are not satisfactory.

- *Encrypted data and models.* The first approach is to train encrypted DNN models over encrypted data. However, due to the large training data and expensive training process in deep learning, cryptographic model training approaches are too expensive to be practical. A recent study on training small-scale neural networks [26] (e.g., just two layers with a maximum of 128 neurons per layer) has shown astonishingly high communication, computation, and storage costs. As a result, cryptographic approaches are often limited to training small models [26], [30] or securely applying trained DNNs for prediction [40], [29], [16].
- *Federated learning.* Another possible solution is to partition the dataset and the learning task into sensitive and non-sensitive partitions and use cloud-client federated learning [19]. The non-sensitive portion, assuming it is much larger than the sensitive one, is exported to and processed by the cloud. Correspondingly, the learning process is partitioned and distributed between the cloud and the client, ensuring that the intermediate information exchanged between the two parties does not breach privacy. Collaborative deep learning framework enhanced with differential privacy [32], [1] may be tweaked into such a partition-based setting. However, reported attacks [14] allow an adversarial collaborator (e.g., a compromised cloud in the cloud-client scenario) to generate images resembling the sensitive classes owned by the victim parties (the trusted data owner).
- *Trusted execution environments.* Hardware-assisted trusted execution environments (TEEs), such as Intel SGX, can also be applied to deep learning in the cloud. The idea is to create a secure enclave in the specific memory area (enclave page cache (EPC)) so that no

other process/thread can access the content in the enclave. Memory pages are also automatically encrypted when they are swapped to the disk. However, recent studies on side-channel attacks [8] make this approach challenging to develop and deploy. Attackers can peek or infer the plaintext content inside the secure enclave via side channels, such as page fault interrupts and cache loading.

Furthermore, to work with GPUs, costly cryptographic approaches have been applied to achieve partial data confidentiality in transferring data between CPU and GPU [35], [27], which does not meet the performance requirement for training large DNNs or with large training data. GPU manufacturers may develop TEEs for GPU¹ However, it's to be tested to determine how secure it is in terms of side-channel attacks.

Researchers have also explored the application of differential privacy (DP) [6] in distributed (federated) learning scenarios [32] or a trusted central training server [1]. However, DP works for the setting of sharing data and models without breaching individual training examples' privacy. It does not meet the need for data and model confidentiality.

Scope and contributions. This study compares two different image disguising approaches: InstaHide [16] and our recently developed DisguisedNets. These image disguising mechanisms protect the training data and also possibly the learned models by casting training data into a confidential transformed space where powerful DNN models can still learn features and patterns distinguishing image classes and leverage the power of GPUs in the cloud. The intuition is twofold. (1) Apply appropriate transformations and data protection mechanisms so that the disguised images cannot be effectively reconstructed and re-link to the original images. (2) Meanwhile, powerful deep learning techniques can still pick up the unique topological and geometric features preserved in the transformed space to distinguish the originally defined classes of images in the transformed space. By doing so, the tie between the original training data and the learned model in the transformed space is broken, which also disables any model-based exploration [3], [28], [34], [10]. In the end, we can approximately preserve the *image distinguishability* for the target classification task while minimizing the *recoverability* of individual images.

There are several unique contributions.

- 1) We have designed two image disguising mechanisms: AES-based (AES) and random-projection-based (RMT) for image-based DNN learning to preserve training data and model confidentiality in outsourced training. The goal is to study and achieve a good balance between the utility of disguised images and the level of confidentiality protection.
- 2) We have carefully analyzed the potential attacks under the outsourced deep learning settings and the resilience of disguising mechanisms to attacks on data and model

confidentiality. So that users can choose the corresponding method under their preferred threat model.

- 3) We have conducted extensive experimental evaluations on public datasets to show the trade-offs of different disguising schemes and related parameter settings between data utility and their resilience to attacks. We also show how the disguising methods work to protect models from model-targeted attacks and

II. RELATED WORK

Sensitive deep learning assets may include training data, models, and online testing data. Protection methods may target the model training phase or the application (i.e., inference) phase when both phases can be exported to the public cloud. However, typically, the computational complexity and the demand on resources of model training is far more than model application. Thus, many studies have focused on the application phase, e.g., a cloud-based model inference service hosting the trained model and making predictions for a user-provided image [12], [29], [17].

In contrast, due to the high computational complexity of training algorithms and the large size of training data, there is no practical cryptographic approach, e.g., homomorphic encryption or secure multi-party computing based approaches for protecting model training. A few recent studies in this direction have shown prohibitively high costs even for a small neural network model [26], [30]. Trusted execution environments (TEE) with masked GPU operations are applied to speed up training [27], [35]. However, no TEE-based deep learning method has addressed severe side-channel attacks [8].

Researchers have looked at protection methods for image training data in the outsourcing context. Noise addition [7], image blurring [24], and morphing [23] are weak as the visible features of the images are still perceivable and understandable. Such transformations do not defeat simple visual re-identification. Figure 1 shows how easy it is to visually re-identify the content in the original images by observing the transformed ones with the mentioned techniques.

A recent method InstaHide [16] applies the idea of mixing up images [41] in a training set with public images with linear combination to obfuscate the content, along with randomized signs of pixel values for further protection. However, it is vulnerable to image reconstruction attacks [2] that need only to know the disguised images and public image sets (i.e., the Level-1 adversarial knowledge, as we will discuss).

Thus, on the one end, existing disguising mechanisms are too weak to protect almost nothing. On the other end, if an encryption mechanism, e.g., homomorphic encryption, or a complex cryptographic protocol, is applied, such linking or reconstruction would be impossible. However, the current cryptographic schemes incur extremely high costs in almost all aspects of computation, communication, and storage. Thus, they are impractical for resource-intensive tasks like training a DNN model. Hardware-assisted approaches, such as TEEs, are still under investigation to ensure the expected security properties. Along with all these possible approaches, we aim

¹Nvidia has announced a GPU TEE in their Hopper architecture.

TABLE I
RELATED WORK ON TRAINING PHASE PROTECTION.

Sample Related Work	Method	Weaknesses	Strengths
Mohassel et al. [26]	Train DNN models over masked or encrypted data with cryptographic protocols.	Involve high communication, computation, and storage costs; Require extensive re-design and custom implementation of the DNN architecture; Involve iterative interactions between the data owner and cloud provider.	Provide semantic security. Model quality is fully preserved.
Tramer et al. [35] and Ng et al. [27]	Use TEEs for confidential CPU operations and masked data for confidential GPU operations	side channel attacks are an unaddressed concern	More efficient than cryptographic protocols, but GPU operations on masked data are still expensive
Abadi et al. [1] and Shokri et al. [33]	Apply differential privacy to randomize intermediate gradients	Aim to share trained models, and thus does not preserve model confidentiality; result in significant drop in model quality; vulnerable to generative adversarial network attacks.	Preserve individuals' privacy.
Fan [7]	Train shallow neural network locally and outsources intermediate representation to the cloud for deeper training.	The intermediate representation of images reveals the visual characteristics of the related images, vulnerable to visual re-identification attacks.	NA
Li et al [24]	Applies differential privacy to hide sensitive pixels in images	Does not hide the global visual characteristic content of images, vulnerable to visual re-identification attacks	NA
Zhang et al. [41] and Huang et al. [16]	Mix-up images from the training set and public domain with random selection and weighing to hide the content of the sensitive image.	Vulnerable to ciphertext-only (Level 1 adversarial knowledge) image re-construction attacks; May expose trained model to wide variety of model-based and membership attacks for the inside-dataset setting.	Fast training; No changes to the training architecture.

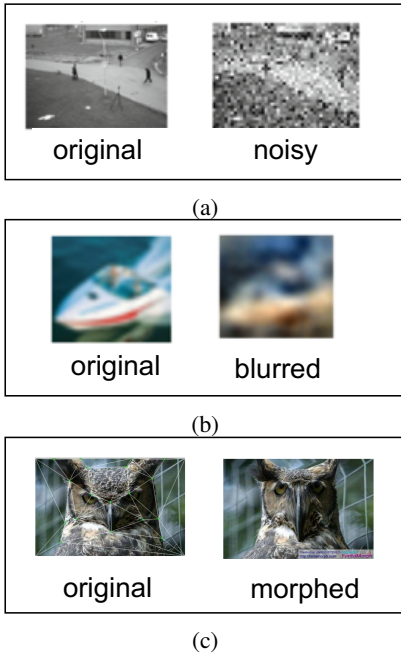


Fig. 1. (a) Differentially private noise addition to images [7]; (b) The reconstructed blurred images in PrivyNet [24]; (c) The morphed images [23]

to explore and develop new image-disguising methods to achieve good balances among costs, data utility, and security guarantees. Table I summarizes the current work in protecting data and models in the training phase.

III. THREAT MODELING

We are concerned with the confidentiality of the sensitive training image data and the DNN models in the *outsourced training phase*. Here, we make some relevant security assumptions for our disguising mechanisms. 1) We consider the cloud provider to be an honest-but-curious adversary, which implies that a curious provider will still honestly deliver desired results to the data owner. However, it may keep a copy of the data and programs it can observe. 2) The adversary can observe the training data, the training process, and the trained models, including the structure of the DNN architecture and parameter settings for training. Thus, they can probe the observed items with methods such as image reconstruction, re-identification, and membership attacks. (3) We do not address evasive attacks and poisoning attacks [3], [28], where adversaries will tamper with the training data, which can be guarded with training data integrity checking. (5) The client infrastructure and communication channels are secure.

Assets to Be Protected. Under a certain protection mechanism, we generalize that the training data D is transformed to $f_{key}(D)$, and the model $M = M(D)$ is changed to $g_{key}(D)$. The attacker might want to know whether an image is likely used to train a model, e.g., the membership inference attack [34], observe training images that may contain sensitive objects, or steal a proprietary training dataset. The attack might also target the model if the model is exposed, e.g., using model-inversion attacks [9] to explore the private information of training data.

Adversarial Prior Knowledge. The adversary may have

two levels of prior knowledge. For each level, we may design a disguising technique.

- *Level-1*: They may know what the model is used for, e.g., the background application, the distribution of the data, e.g., face images, and the type of disguising technique used, but do not know the disguising parameter setting for a specific dataset that serves as the secret key to the protocol.
- *Level-2*: In addition to Level-1 knowledge, they may try to obtain pairs of images and their disguised versions via other attacking channels (not including the ones they are targeting). They hope to use these known pairs to explore various image reconstruction attacks.

Potential Attacks. Recent studies have shown that attackers can explore training/testing examples and models, for example, to find adversarial examples misleading the prediction of deep neural networks [3], [28]. Such attacks depend on adversaries' clear understanding of the original image data and the ability to access the developed models freely. Outsourced learning without protection makes these attacks easier to deploy.

With a protection mechanism on data and models, we consider a fundamental attack: *training image re-identification* that aims at linking the protected images to identifiable original images. We introduce a model-based re-identification test – DNN examiner, which uses a model trained on the original data to tell whether a protected image is re-identifiable. Note that some related methods [7], [24] are not resilient to human visual re-identification, which does not protect confidentiality. Since the image disguising mechanisms break the link between the original training data and the learned models (in the transformed space), the existing model-oriented attacks do not work anymore without successfully breaking the disguising mechanism. Attackers thus depend on *reconstruction attacks*: reverse the disguising mechanism to approximately reconstruct the original images and then try to re-identify the reconstructed images. We use the DNN examiner approach to evaluate how successful a reconstruction attack works in our experiments.

IV. DISGUISEDNETS – A NOVEL IMAGE DISGUIISING MECHANISM FOR OUTSOURCED DEEP LEARNING

In the following, we will introduce an image disguising framework that incorporates pixel-block partitioning, random block permutation, and block-wise transformations of images along with noise additions. The premise is that after the dramatic transformation, it is difficult to link the disguised images to the original images, while, unlike pure encryption schemes, it still preserves some essential patterns for distinguishing between classes of images that allow DNN learning methods to capture. This amalgam of multiple transformations provides a sufficiently large parameter space so that the attacks are computationally intractable (Section V) under the Level-1 prior knowledge.

Figure 2 depicts the DisguisedNets framework. A data owner disguises her private images before outsourcing them to the cloud for DNN learning. She can either fully outsource the entire image datasets and the learning procedure to the

cloud or selectively retain sensitive images in the cloud-client partitioning setting. She transforms all of her images using one secure transformation key secret to her. Note that this transformation should be at a reasonable cost, practical for a client's infrastructure to process.

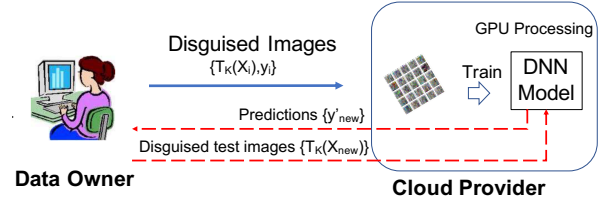


Fig. 2. Image disguising framework for DNN learning.

Specifically, assume the data owner owns a set of images for training, notated as pairs $D = \{(X_i, y_i)\}$, where X_i is the image pixel matrix ($l \times m$ and $3 \times l \times m$ for grayscale and RGB images respectively) and y_i the corresponding label. We formally define a disguising mechanism as follows. Let the disguising mechanism be a transformation T_K , where K is the secret key that depends on the selected perturbation techniques. By applying image disguising, the training data is transformed to $\{(T(X_i), c_i)\}$ with c_i mapped to $0, 1, \dots$ randomly representing the classes y_i . The original model is a function $\mathbf{M}(D)$, which is learned with a DNN learning method \mathcal{G} : $\mathcal{G}(D) \rightarrow \mathbf{M}(D)$. The image disguising mechanism enables the same learning method \mathcal{G} to be applied to the transformed data directly without any modification: $\mathcal{G}(D_T) \rightarrow \mathbf{M}_T(D_T)$. For any new data X_{new} , the model application (or inference) is defined as $\mathbf{M}_T(T(X_{new}))$, i.e., the new data transformed with the same key. To make such a transformation method practical for modeling, i.e., a model trained with transformed data still working satisfactorily, a user may expect the error of modeling is not far away from the original model's. Thus, a utility-preserving mechanism should have

$$|Err(\mathbf{M}_T) - Err(\mathbf{M})| < \delta$$

where δ is the level of model quality degradation acceptable to the user. While for a specific DNN modeling method and a specific dataset, it's difficult to theoretically justify what this gap will be, one can always directly evaluate the model quality to check whether it is acceptable for the application. We have empirically evaluated the δ levels for different mechanisms, datasets, and a few popular DNN modeling architectures in experiments.

A. Pixel-Block Partitioning and Block-based Random Permutation

In this section, we present one way of image transformation: image block permutation, that will be combined with other mechanisms later.

An image $X_{l \times m}$ is first partitioned into t blocks of uniform size $r \times s$. If we label the blocks sequentially as $v = \langle 1, 2, 3, 4, \dots, t \rangle$, a pseudorandom permutation of the image, $T_\pi(X)$, shuffles the blocks and reassemble the corresponding

image accordingly. Block-based permutation preserves the in-block information and the relative positions of related blocks. Thus, we understand it preserves a great amount of information for effective modeling. However, while the permutation may break the global patterns of the images and achieve good visual privacy already, the between-block characteristics such as boundaries, color, content shape, and texture of the original neighboring blocks may provide clues for adversaries to recover the original image – imagine the jigsaw puzzle! For large t , such attacks can be time-consuming due to the vague similarity between block boundaries. However, with the prior knowledge: a pair of original image and its block-permuted image, it’s not difficult to solve such a jigsaw puzzle. Thus, we use this as an auxiliary step enhancing other steps in the disguising framework.

B. Pixel-Block Transformations

Next, we establish pixel-block-level protection mechanisms that aim to preserve the data utility for DNN modeling and further increase the resilience to attacks. We consider two candidate mechanisms: random projection and encryption schemes, and discuss their characteristics. Specifically, when an image is partitioned into t pixel blocks for random permutation, we get a list of t parameters $\{K_i, i = 1 \dots t\}$, one for the pixel-block at the same position across the whole dataset. We name the specific position of the pixel block in the image *the pixel-block position*. The list of parameters acts as a secret key and will be shared, together with the permutation key, by each image in the dataset. The purpose of this setting is to maximize the preservation of distinguishable patterns between image classes – i.e., a pair of similar image patterns (blocks) can still be transformed to another pair of (likely) similar ones after applying the disguising mechanism.

Randomized Multidimensional Transformation (RMT).

For an image represented as a pixel matrix X , a general linear transformation can be defined as $G(X) = R(X + \Delta)$, where $R_{m \times m}$ is a random orthogonal matrix generated following the Haar distribution [11], or a random invertible matrix, e.g., a random projection matrix [39], and Δ is an optional noise matrix. We call this method the randomized multidimensional transformation. When an image is partitioned into t blocks for random permutation, we prepare a list of random matrices $\{R_i, i = 1..t\}$, one for each image-block position and share this list for each image. Such transformation is known to preserve (or approximately preserve by random projection) the Euclidean distance between columns of the matrix X . For real application, we may arrange the pixel blocks accordingly to form the column of X . For example, a 4×4 pixel matrix can be partitioned into 4 2×2 block to preserve the smaller block-level similarity with RMT. Figure 3 shows the effects of RMT on MNIST and CIFAR-10 datasets.

AES Block Transformation (AES). The existing AES encryption schemes typically use 128-bit encryption keys, which encode every 16-byte data block sequentially. If we use AES for pixel-block encryption, assuming each pixel is stored in one byte, 16 original pixels are mapped to 16 encrypted

Algorithm 1 DN_RMT (X, t, Key)

Require: X : image of size $l \times m$; t : number of blocks; $\text{Key} = \{\text{permutation_key}, \text{transformation_matrices}, \text{noise_level} \in [0, N]\}$

- 1: $r, s \leftarrow$ compute image block size with $l \times m$ and t ;
- 2: Partition image $X_{l \times m}$ into blocks X_1, X_2, \dots, X_t ;
- 3: Shuffle the image blocks pseudorandomly with permutation_key
- 4: **for** each block $i, i = 1 \dots t$ **do**
- 5: $\Delta_i \leftarrow$ Generate random matrix with elements from the uniform distribution in $[0, N]$;
- 6: use the transformation matrix at the position i : R_i ;
- 7: $Y_i \leftarrow R_i(X_i + \Delta_i)$;
- 8: **end for**
- 9: Re-assemble $\{Y_i\}$ to make the transformed image Y and return Y ;

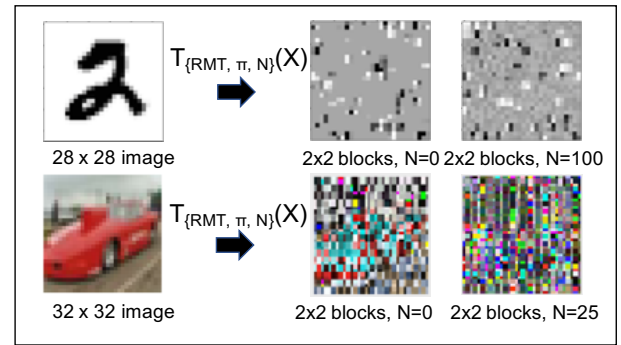


Fig. 3. Block-wise RMT+Noise on MNIST and CIFAR-10 images.

bytes (pixels), and a whole pixel block is encoded to 16-byte units. Putting all encrypted pixel blocks together, we get a disguised image. For clear presentation, when we talk about AES encryption block, i.e., 16 bytes for a 128-bit encryption key, we use the 16-byte “encryption unit”, which are different from “pixel blocks” we have been using previously in our image disguising framework. Figure 4 shows some example AES transformations on images.

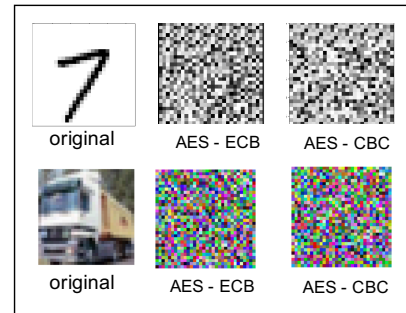


Fig. 4. Pixel-block based AES encryption of MNIST and CIFAR-10 images.

We consider two AES modes in our design. (1) We observe that with the AES Cipher Block Chaining (CBC) mode, any pixel-level change in the pixel block between two images

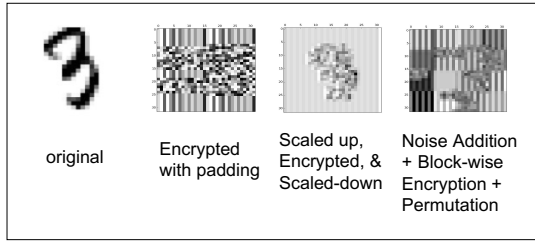


Fig. 5. AES-ECB encryption of MNIST image with different strategy.

will result in different encoding results for most 16-byte AES blocks in this pixel block position, making it not ideal for our purpose. (2) Then, we turn to the AES Electronic Code Book (ECB) mode that can be considered as a fixed mapping function between 16-byte original data to 16-byte encrypted data. Different from CBC, the neighboring 16-byte blocks do not affect the encoding of the current block. This matches our requirement of data utility preservation, e.g., to preserve the block-level distinguishable patterns after the transformation.

To preserve more information, based on the intuition of smaller blocks preserving more inter-pixel-block information, we can also use unit sizes smaller than the regular size, which is 16 pixels for 128-bit ECB. The method is to scale up the image first, e.g., from 32x32 to 256x256 (where each pixel is duplicated eight times), then encrypt it by 16-pixel units, and finally scale down to the size 32x32. Please refer to Figure 5 for the detailed example. It's equivalent to encrypt 2-pixel units in the original 32x32 images. We found by reducing the block size, the model quality can be improved with the cost of lower attack resilience to the Level-2-knowledge-based attack.

Algorithm 2 DN_AES (X, k, t, Key)

Require: X : image of size $l \times m$; k : scale-up factor; t : number of blocks; $\text{Key} = \{\text{permutation_key}, \text{AES_keys}, p = \text{probability of salt-pepper noise}\}$

Ensure: the selection of k and t results in image blocks that can be further partitioned to 4x4 pixel patches;

- 1: $X_{lk \times mk} \leftarrow$ scaled up the image;
 - 2: $r, s \leftarrow$ compute image block size with $lk \times mk$ and t ;
 - 3: Partition image $X_{l \times m}$ into blocks X_1, X_2, \dots, X_t ;
 - 4: Shuffle the image blocks pseudorandomly with permutation_key;
 - 5: **for** each block $i, i = 1 \dots t$ **do**
 - 6: $Y_i \leftarrow$ for each pixel in block X_i , with the probability p , it's randomly turned to white or black pixel (salt-pepper noise);
 - 7: $E(Y_i) \leftarrow$ with the AES CBC mode, every 16-byte segment (4x4 pixel patch) is encrypted to 16 bytes of AES digest with AES_key_i;
 - 8: **end for**
 - 9: re-assemble image blocks $E(Y_i)$ and return $E(Y)$
-

C. Complexity Analysis

The additional costs of the disguising methods consist of the encoding cost and the possible additional learning cost, i.e., it may take more rounds to converge. We leave the second part to the experimental evaluation and analyze the encoding cost here.

For an image partitioned into t blocks with each block $l \times m$, the RMT transformation involves t matrix-matrix multiplications and matrix additions. As the numbers t, m , and l are all small, the cost of RMT per image is low: $O(tlm^2)$. For an image of $l \times m$ with a scale-up constant of s , the AES-128 encryption cost is $l \times m \times s/16$ times of AES encryption. Our experimental evaluation shows that per image cost is less than 10 ms and can be comfortably done by any PC or mobile phone.

D. Model Protection via Image Disguising

Note that the models trained with disguised data work only on disguised data. We show this property also protects models from existing model-targeted attacks. So far, we have seen model-inversion attacks [9], [42], membership-inference attacks [34], [15], and model-extraction attacks [36], [18].

Model-extraction attacks assume the attacker can freely access the model, e.g., via a cloud-based prediction API. With such a service, the attacker can try various images to collect their outputs and then use the input-output pairs to reconstruct the model. Our threat model assumes the attacker can copy or save the trained model for analysis. Thus, the attacker does not need to perform model-extraction attacks. As the models only work on disguised test images, without the secret disguising key, they are useless to the attacker.

Membership-inference attack aims to estimate the possibility of a target example belonging to the training data of a model. To perform such an attack, the attacker must first apply the disguising method (with the secret key) to the target example so that the model can be used. This step effectively blocks the attack or at least significantly increases the difficulty. To successfully conduct the MIA attack on the disguised model, the attacker may need to manipulate an authorized user to transform the example and intercept the transformed one, which corresponds to the mentioned Level-2 knowledge. Thus, the disguising mechanism establishes an effective defensive line.

Model-inversion attack uses a learning procedure, e.g., a GAN method [42], to progressively adjust randomly generated or seed images from similar domains towards most likely training examples. When applied to the models trained on disguised data, the model-inversion attack recovers only the disguised training data, not the original data. Again, the disguising mechanism builds a defensive line on this attack. We will show how the RMT mechanism works against model-inversion attacks in experiments.

V. ATTACK ANALYSIS

This section aims to analyze the possible threats to the proposed disguising mechanisms and clarify the applica-

ble settings. With Level-1 adversarial knowledge, Disguised-Nets mechanisms provide strong confidentiality protection, as shown in the discussion of “brute-force attacks”. In contrast, other related methods are still struggling with visual re-identification by human eyes [24], [7] or disguised-image-based reconstruction attacks [2]. We also analyze more sophisticated reconstruction attacks that depend on Level-2 adversarial knowledge.

A. Level-1 Adversarial Knowledge and Attacks

Recall that Level-1 knowledge includes knowing the disguised images and possibly the model domain, i.e., the types of images and the background application. It is clear that with only Level-1 knowledge, the brute-force attack on AES schemes is not possible, and thus we focus on the scheme using multidimensional projection.

Visual Re-identification. The first simple attack is to visually identify images by human attackers. We have shown that simple methods like noise addition, morphing, and shallow-network-based transformation are not resilient to this attack. However, many other attacks may use re-identification as the last step, i.e., reconstruction attacks. It’s inefficient for human evaluators to check each image to determine the protection level of an image disguising mechanism. Thus, we propose the *DNN examiner* approach for evaluation purposes: let a DNN trained on the original datasets to perform the visual re-identification task for human evaluators. We will use DNN examiners in experiments.

Brute-Force. The brute-force attack method for image reconstruction is to enumerate each possible parameter setting of the disguising mechanism and then check the recovered result with re-identification. As AES encryption is already resilient to the brute-force attack, we examine the RMT method only. Let’s start with a block-level transformation for any image block i with RMT. With $X'_i = X_i R_i$, the adversary knows only X'_i . In the brute force attack, the number of possible X_i is determined by the number of possible R_i matrices. We show that the number of possible R_i (even limited to orthogonal ones) can be exponentially large for given parameters.

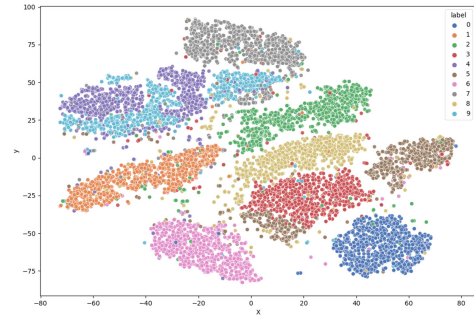
Proposition 1. For values encoded in h -bit finite field, there are $O(2^{hm})$ candidate orthogonal matrices $R_{m \times m}$.

The proof is based on the theory of orthogonal matrix group [13], the detail of which is skipped here. With a typical setting in our experiments, e.g., $h = 8$ and $m = 28$ for the MNIST dataset, the overall complexity is $O(2^{224})$, which is more than sufficient to protect from computationally-bounded attackers. Combined with the random permutation of blocks, the attack complexity is even higher. Thus, a brute-force attack is generally impractical for the proposed methods.

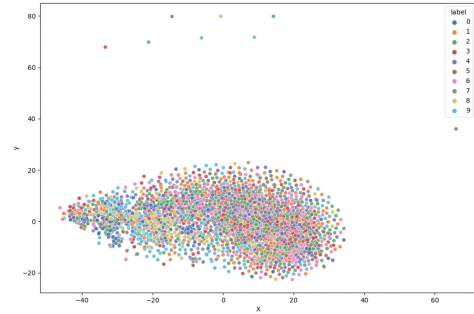
Clustering Attack. Carlini et al. [2] utilized a clustering method to attack InstaHide [16] disguised images. InstaHide uses the random mix-up method to generate disguised images. Depending on the random weight distribution, some disguised images might be dominated by the same image, which likely forms a cluster of images that can be used to de-mask and

de-noise. As InstaHide disguised images are essentially linear combinations of plaintext images, the attack result can be visually re-identified.

Important questions are whether our disguising methods can generate images with clustering structures and whether such clusters can be used to break our disguising methods. To answer these questions, we visualize the disguised training data with t-SNE [38] to understand the existence of clustering structure in the Euclidean-distance space. Figure 6 shows that RMT might preserve the clustering structures for some datasets: for simpler datasets like MNIST and FASHION, the clustering structure is well preserved, while others do not. In contrast, AES does not preserve any clustering structure, as shown in Figure 7. While AES not preserving clustering structures to leave less information to attackers, it also affects data utility and leads to lower-quality models, as we will show in experiments.



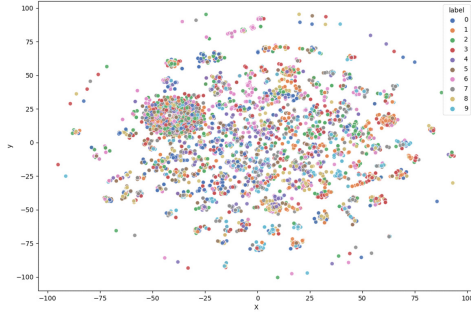
(a) RMT on MNIST



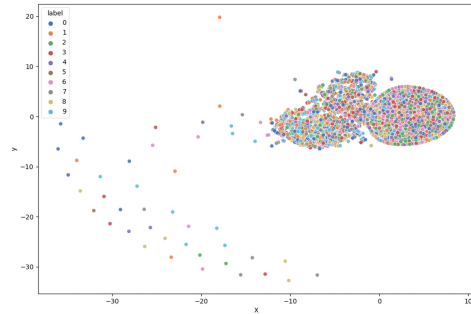
(b) RMT on CIFAR10

Fig. 6. t-SNE visualization of RMT disguised datasets (4x4 blocks). Colors represent different labels. A dense area covered with one color means that the clustering structure matches the label distribution well for the specific subset.

Next, can such preserved clustering structures be used for attacks? An attack on InstaHide [2] has used image clusters to de-noise and de-mask, as InstaHide uses the random weights mix-up mechanism. However, unlike InstaHide, clustering structures of RMT-disguised images do not help attackers identify original images. However, it may help attackers infer additional information with other domain knowledge. For example, if the original training samples’ distribution (and



(a) AES on MNIST



(b) AES on CIFAR10

Fig. 7. t-SNE visualization of AES disguised datasets (4x4 blocks)

thus the clustering structure) is also known, it may allow the attacker to identify the mapping between a specific cluster of original images and a cluster of disguised images. As distances between samples are not preserved, it’s still difficult to figure out the sample-to-sample mapping.

B. Level-2 Adversarial Knowledge and Attacks

We move one more step further to study the more challenging issue: what if a powerful adversary can obtain additional knowledge: pairs of original images and their transformed ones? This assumption corresponds to the chosen-plaintext attack in cryptographic analysis [20]. This study helps us understand when we should not use a proposed disguising method. Below, we focus on the *codebook attack* on the AES-ECB-based disguising mechanism and the regression-based attack on the RMT mechanism.

Codebook Attack. The assumption is that the adversary is knowledgeable of the encryption procedure described previously but does not have the encryption/decryption key. Since the AES ECB method is deterministic, the basic attack is to build a mapping (i.e., the codebook) between the plaintext unit (e.g., the 16-byte pixel block) and its encrypted counterpart (e.g., the 16-byte AES cipher block). By processing the known image pairs, the adversary constructs a codebook as a dictionary mapping 16-byte pixel blocks to encrypted 16-byte blocks. Since different images, especially those in the

same class, might share some 16-byte pixel blocks, some 16-byte encrypted blocks in the targeted images are likely already in the codebook, which will be used to recover the original blocks. For encrypted pixel blocks not present in the codebook, the adversary may use a fixed pattern, e.g., all zero values or most likely values to pad. By repeating this procedure for each 16-byte block, the adversary can recover some parts of the image, which can be further re-identified via human eyes or models at the adversary’s hands.

Possible mitigation methods. Let the *hit rate* be defined as the probability that an encrypted pixel block can find a match in the codebook. This attack can become less effective if we add salt-and-pepper noises to the original images before encoding. This step will reduce the hit rate significantly and make the mapping non-unique: the same 16-byte pixel block can be mapped to different ciphertexts. We will evaluate the *success rate* of this attack in experiments, using the accuracy that the DNN examiner trained with the original image space correctly classifies the reconstructed images.

Projection Matrix Estimation Attack. Note that noise addition can easily defend the RMT method from the codebook attack, which is already a part of the RMT method. However, if the adversary has obtained enough original and transformed image pairs, there is a possibility that the transformation matrix might be estimated with linear regression. Specifically, a noise-added block-wise transformation, e.g., $Y_i = R_i(X_i + \Delta_i)$, where Δ_i is a random noise matrix, re-generated for each image block X_i , and drawn uniformly at random from $[0, N]$ where N is the tunable noise level. With enough known pairs of (X_i, Y_i) , the regression method can be applied to estimate R_i . Generally, the more known pairs, the more precise the estimation can be. However, it’s unclear how the noise level affects the effectiveness of estimation and how we can achieve a good balance between data utility and attack resilience. We will examine the regression-based attacks in the experiments.

Note that the recently proposed InstaHide [16] method also somewhat matches this definition of image disguising. It also requires learning from the disguised examples $\{(T(X_i), c_i)\}$. However, the learned model M_T is still applicable to the original test data, i.e., the application phase uses $M_T(X_{new})$. They also show that the performance of $M_T(X_{new})$ is very close to $M(X_{news})$, which implies $M_T \approx M$. leads to serious problems, such as the impossibility result and a clustering-based attack, as Carlini et al. [2] show. In contrast, our proposed methods require strictly $M_T(T(X_{new}))$ in the application phase, which eliminates the possible information leakage targeting the models and the clustering of disguised training images.

VI. EXPERIMENTS

The experiments have three goals. (1) The proposed DisguisedNets mechanisms involve parameter settings, which may affect data utility. (2) While the proposed methods are resilient to attacks under Level-1 knowledge, we need to understand the intrinsic trade-offs between data utility and the methods

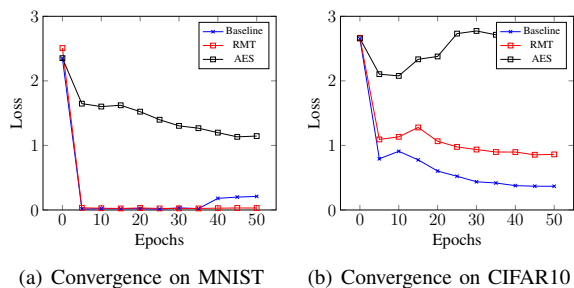


Fig. 8. Convergence speed on disguised images. Baseline: models trained on original datasets.

enhancing the resilience to Level-2 attacks. (3) As we have discussed, the proposed methods have unique benefits in defending model-based attacks, which we will demonstrate in experiments.

Datasets. We use four prevalent DNN benchmarking datasets: MNIST, FASHION, CIFAR10, and LFW [22] for experiments. MNIST (handwritten digits) and FASHION (fashion items) are gray-scale 28×28 -pixel images with ten classes. CIFAR10 has 60 thousand 32×32 color images distributed into ten classes. LFW is a labeled face database. It is relatively small, with only a few thousand samples. We used five folds of random sampling to estimate the standard deviation of modeling results, which are also used for later experiments.

Table II summarizes the datasets, the techniques used to train the base models, and their baseline model accuracy on the original image data. All the models are implemented with PyTorch.

TABLE II
DATASETS AND BASELINE ACCURACY. TR: TRAINING, TE: TESTING

Datasets	Records	ImageSize	Network	BaselineAccuracy
MNIST	(60K Tr, 10K Te)	$\{28 \times 28\}$	AlexNet	$96.7 \pm 0.2\%$
FASHION	(60K Tr, 10K Te)	$\{28 \times 28\}$	AlexNet	$88.7 \pm 0.3\%$
CIFAR-10	(50K Tr, 10K Te)	$\{32 \times 32\}$	ResNet-18	$93.4 \pm 0.2\%$
LFW	(1164 Tr, 292 Te)	$\{60 \times 48\}$	ResNet-18	$94.3 \pm 2.0\%$

A. Parameter Settings for Level-1 Attacks

Since all the proposed methods are resilient to Level-1 attacks, we focus on the utility preservation aspect in this section.

Costs. The disguised images are used directly with the existing DNN training algorithms without any modification to the algorithm or data. We have briefly analyzed the per image disguising cost in Section IV-C, which can be comfortably handled by a mobile phone. Another question is whether the disguised images will extend the training time. Fig 8 shows the evaluation of convergence speed on MNIST and CIFAR10 for the three methods: the baseline, RMT, and AES. The baseline refers to the models reported in Table II. Both RMT and AES run with the basic setting of 4×4 blocks. All of the methods converge with 50 epochs, but AES appears more unstable on CIFAR10.

RMT Mechanism. We look at the effects of block size and noise levels on models trained on images transformed with

RMT methods. For easier presentation, we convert block size into the number of blocks: 1 block on the x-axis means the image is not split into blocks; while 196 blocks means 196 2×2 blocks for 32×32 images (CIFAR10) or padded 28×28 images (MNIST and FASHION), and 196 4×3 blocks for padded 60×48 images (LFW). Thus, a smaller block size results in a larger number of blocks after partitioning, as the image size is fixed. If more than one block is generated in partitioning, we also apply a secret block-wise permutation. Figure 9 (a) shows that the model quality is slightly decreased with smaller block sizes (more blocks per image). Overall, the model quality is well preserved, only 2-3% worse than the baseline. It's also understandable that the simpler images, MNIST and FASHION, are more resilient to noise addition and more sophisticated ones are sensitive to noise as shown in Figure 9 (b).

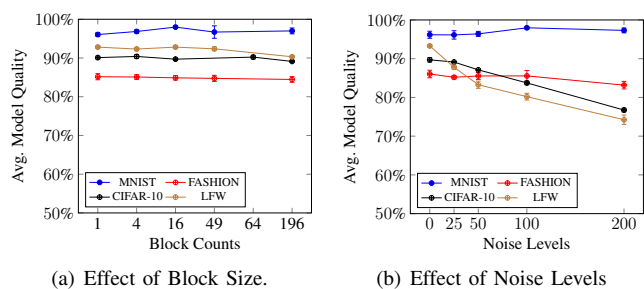


Fig. 9. Effects of block size and varying noise levels on model quality for RMT-disguised images.

AES Mechanism. We have done experiments to understand the effect of block-size setting for the AES ECB based block protection. We use “pixel blocks” for partitioning and permutation, and “units” for AES encryption units. A pixel block typically contains more than one unit. Recall that AES uses 16 bytes as the encryption unit if 128-bit encryption is used. Our partitioning schemes follow this restriction of unit size to make sure that each block has integer times of 16 bytes. Figure 10 shows different block size settings from 1 block (e.g., 32×32 per block for 32×32 images) to 64 blocks (e.g., 4×4 pixels per block for 32×32 images).

We tested two schemes: no scaling vs scaling. The no-scaling scheme uses the block size ≥ 16 bytes, while scaling can use even smaller block sizes. Specifically, when we use a block size ≤ 16 , e.g., 2×2 blocks, the scaling up factors are determined for the x and y axes, corresponding, e.g., the scaling factor for x -axis is 2 and also 2 for y -axis for 2×2 blocks, so that we can partition the scaled image with 4×4 blocks. Figure 10 (a) shows that the model quality can be affected by the no-scaling scheme. For some datasets, e.g., CIFAR10 and LFW, the model quality is too low to be used. Figure 10 (b) shows that the model quality is boosted to the level comparable to the RMT's results for MNIST and FASHION, while the other two still stay at unusable levels. The possible reason is that the colored (multi-channel) images contain more noisy image blocks, which changes significantly after the AES transformation. In summary, different from the

RMT scheme, the AES scheme may only work for some datasets.

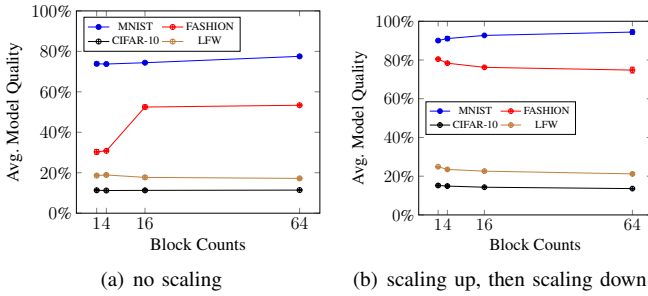


Fig. 10. Effect of block size on model quality for AES-ECB-disguised images

B. Resilience on Level-2 Attacks: AES Scheme

With the known additional knowledge, i.e., pairs of original and disguised images, the disguising mechanisms might be under the reconstruction attack, and attackers can visually check the reconstructed images to re-identify the features of original images.

To effectively evaluate the re-identification step, we use a DNN trained on the original image data to simulate the attacker in the visual re-identification process. The intuition is that if any features in the disguised (or reconstructed) images can be detected visually by the adversary, it can be used to link the disguised (or reconstructed) images to the original images. Such linking is often probabilistic, and we can use the linking *success rate* (i.e., the accuracy of prediction) to gauge the threat level of attacks. As DNNs perform comparably well as human experts do in the image-based classification tasks [21], we believe such “DNN examiners” can satisfactorily simulate the attacker.

We train DNN examiners with the original training data using the same DNN architectures detailed in Table II. We then apply the DNN examiners to see whether the reconstructed images can be correctly classified to their original labels. To minimize the impact of DNN architecture and different baseline accuracy, we define the *attack success rate* as

$$\frac{\text{accuracy of DNN examiner on disguised/reconstructed images}}{\text{accuracy of DNN examiner on original images}} \times 100.$$

1) Resilience to Codebook Attack for the AES ECB method:

Assume the attacker knows m pairs of original images and their ECB encrypted ones, and also other information such as their pixel-block sizes. The codebook attack uses the known pairs to construct a mapping between the known plaintext 16-byte pixels (or a reduced number of pixels if the scaling up/down method is used to preserve more utility) and the corresponding encrypted 16-byte pixels. The attacker might be able to use the codebook to partially recover the original pixel blocks of a disguised image (with random pixel patches for unrecognized blocks). We use the DNN examiner to examine the quality of reconstructed images.

As MNIST and Fashion perform reasonably well with the AES scheme (Figure 10 compared to the other two, we pick

only the MNIST data for clear presentation – the Fashion data has a similar pattern. Figure 11 compares the attack results on 16-pixel encryption units (subfigure (a)) and 2-pixel encryption units with scaling (subfigure (b)). The attacker’s known pairs are selected randomly from the training data, while the targeted images are selected from the testing data. 16-pixel encryption unit gives a one-to-one mapping between the original pixel units and the encrypted ones. We observed hit rates are quite low (lower than 10%), but success rates are increasing steadily due to the increased codebook size. Overall, attackers will need a large number of pairs to achieve a good success rate. 2-pixel encryption unit may create a one-to-many mapping between original pixel units and the encrypted ones, due to the scale up/down processes. We used the Python library function for image scaling. With the scaling process, we observed that hit rates initially increase to around 10% and then drop to 2-3%. However, the success rate quickly reaches the plateau – around 50% with only 20 image pairs. Therefore, no-scaling method is more resilient to attacks – both the hit rates and success rates grow slowly and knowing the whole training data does not help improve the success rates much. In contrast, the scaling method can help gain better model quality. However, it might be vulnerable to Level-2 attacks. There seems an abrupt trade-off the user may have to make.

Aiming at achieving a better balance of utility and attack resilience for the setting of the 2-pixel encryption unit, we found that it’s possible to defend from the codebook attack by adding “salt-and-pepper” noises to the original images. The AES encrypted pixel block changes dramatically when any of the original pixel changes, which helps reduce the attack success rate. Figure 12 shows by adding a small amount of noise, e.g., 2-3%, the attack success rate drops by 10%, while the model quality is not significant damaged. Certainly, the level of noise should be carefully chosen to avoid destroying the data utility: an increase of noise intensity to 4% will dramatically degrade the model quality as Figure 12 (b) shows.

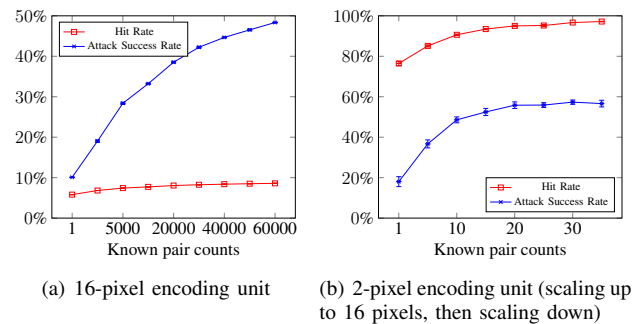


Fig. 11. Codebook attack on MNIST dataset with varying number of known pairs.

C. Resilience on Level-2 Attacks: RMT Scheme

We study how known pairs can be effectively used to attack the RMT method. Again, we assume a stronger attack scenario: the attacker already knows the pixel-block size and the permutation pattern. By known only one pair of

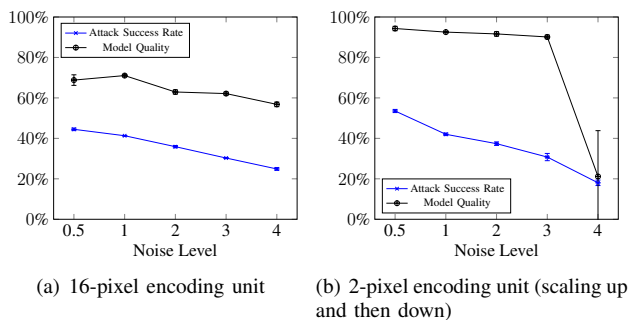
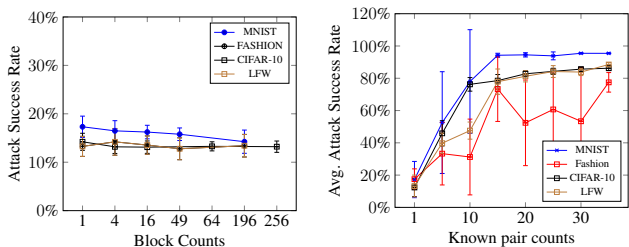


Fig. 12. Protecting AES-based disguising with noise addition (MNIST data)



(a) direct re-identification attack is not effective (b) Regression attacks on the noise-added RMT disguising method can be effective with enough known pairs. Images with block-size 7×7 and noise level $u = 100$

Fig. 13. Attacks on RMT-disguised images.

images, RMT without noise addition can be easily broken – the block-wise transformation parameters $\{R_i, i = 1..m\}$ can be straightforwardly recovered. Different from the “salt-and-pepper” noise for selected pixels in the enhanced AES scheme, we generate a noise value for each pixel and add it to the original pixel value before applying the projection, i.e., $Y_i = (X_i + \Delta_i)R_i$, where the noise Δ_i is drawn from the uniform distribution $U(0, u)$. With noise addition, the known attack method is to use linear regression to estimate the parameters $\{R_i\}$, the accuracy of which is affected by the noise intensity (i.e., the variance of noise) and the number of available pairs.

Figure 13 (a) shows that direct re-identification (with Level-1 attack) is generally not effective at all. However, Figure 13 (b) shows that the regression attack is surprisingly effective on all datasets. With a small number of known image pairs, the attack can achieve surprisingly high success rates. Thus, it’s not safe to use the RMT scheme when Level-2 attack knowledge is possibly available.

D. Use Image Disguising to Protect Models

Exposing models may have high risks, as shown in model-inversion attacks, membership-inference attacks, and model-extraction attacks. This experiment shows that image-disguising methods can work effectively against such model-targeted attacks. We take model-inversion (MI) attacks, for example, which try to recover training data from the exposed model.

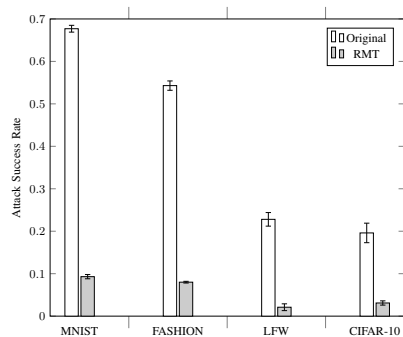


Fig. 14. RMT protects models from model-inversion attacks. Original: the MI attack applied to the original non-protected model to recover images. RMT: the MI attack applied to the model trained on RMT-disguised training data.

TABLE III
BEST RESULT UNDER LEVEL-1 ASSUMPTION

Datasets	No Disguise	RMT Disguise	AES Disguise
MNIST	96.7 \pm 0.2%	96.6 \pm 0.4 %	91.6 \pm 1.1%
FASHION	88.7 \pm 0.3%	85.1 \pm 0.6 %	68.9 \pm 1.4%
CIFAR-10	93.4 \pm 0.2%	89.3 \pm 0.1%	11.3 \pm 1.7 %
LFW	94.3 \pm 2.0	92.6 \pm 2.3%	17.2 \pm 0.4%

The experiment exposes the models trained on RMT disguised images for a recent model-inversion (MI) attack [42] that has shown good performance in recovering training data. Specifically, we used a 4×4 block without noise addition for RMT to generate disguised data and models. We then apply the MI attack to generate 2000 images for each dataset (200 for each class). To compare the performance of the MI attack, we use the DNN examiners trained on the original datasets to recognize the recovered images. Figure 14 shows that models trained on RMT-disguised data are very resilient to the MI attack. Indeed, the MI attack recovers the RMT-disguised training data, which are different from the original images and thus still unrecognizable. The results are consistently worse than the DNN examiners applied to the RMT-disguised training data directly (Figure 13 a).

E. Discussion

Based on the experimental results, we have the following observations.

- With Level-1 adversarial knowledge, the RMT mechanism preserves good data utility for most datasets. In contrast, the AES scheme only keeps data utility for some datasets. Table III summarizes the best result under the Level-1 adversarial knowledge assumption.
- With Level-2 adversaries, the RMT mechanism should not be used as the attack success rate will be high. The AES scheme with a small encryption unit and small (e.g., 2%) noise addition is resilient to the codebook attack and still preserves model quality for some datasets. Table IV summarizes the best results for AES. As the AES scheme does not work on CIFAR10 and LFW, so far, we haven’t discovered satisfactory utility-preserving disguising methods against Level-2 adversaries.

TABLE IV
AES BEST RESULT UNDER LEVEL-2 ASSUMPTION: ENCRYPTION UNIT
2X1 (WITH SCALING), NOISE LEVEL 2%.

Datasets	No Disguise	Model Accuracy	Attack Success Rate
MNIST	96.7 \pm 0.2%	90.14 \pm 1.1%	30.76 \pm 0.87%
FASHION	88.7 \pm 0.3%	73.08 \pm 0.86%	23.51 \pm 0.27%

- Finally, if only Level-1 adversaries are expected, RMT can also be used to effectively protect from model-targeted attacks, as the models trained on RMT disguised data can only be applied to disguised data.

VII. CONCLUSION

Outsourcing large image datasets to the cloud for deep learning has been an economical and popular option, but it also raises concerns about data and model confidentiality. The existing solutions are either too expensive to be practical, vulnerable to different model-based adversarial attacks, or ineffective in protecting the image content. By focusing on the training image reconstruction and re-identification attacks, we propose image disguising mechanisms that efficiently thwart the attacks and preserve model quality. The combination of random image-block permutation and block-wise AES encryption or multidimensional transformation (RMT) does not require any changes to the existing DNN modeling architectures. Experimental results show that the RMT method can preserve the model quality and provide sufficient attack resilience under Level-1 adversarial knowledge – adversaries knowing only the disguised images and the domain information. The AES method improves the attack resilience against Level-2 adversaries who manage to obtain pairs of original images and disguised ones. However, the AEs method may seriously damage some datasets’ utility. We also show that the disguising methods can protect the trained models from model-targeted attacks. The future work will be focused on new image disguising mechanisms that can more efficiently preserve utility with stronger security guarantees. We will also extend the related research to non-image data.

REFERENCES

- [1] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, CCS '16*, pages 308–318, New York, NY, USA, 2016. ACM.
- [2] N. Carlini, S. Deng, S. Garg, S. Jha, S. Mahloujifar, M. Mahmood, S. Song, A. Thakurta, and F. Tramèr. Is private learning possible with instance encoding? In *IEEE Symposium on Security and Privacy (S&P)*, 2021.
- [3] A. Chakraborty, M. Alam, V. Dey, A. Chattopadhyay, and D. Mukhopadhyay. Adversarial attacks and defences: A survey. *CoRR*, abs/1810.00069, 2018.
- [4] A. Chen. Gcreep: Google engineer stalked teens, spied on chats. *Gawker*, <http://gawker.com/5637234/>, 2010.
- [5] A. J. Duncan, S. Creese, and M. Goldsmith. Insider attacks in cloud computing. In *2012 IEEE 11th International Conference on Trust, Security and Privacy in Computing and Communications*, 2012.
- [6] C. Dwork. Differential privacy. In *International Colloquium on Automata, Languages and Programming*, pages 1–12. Springer, 2006.
- [7] L. Fan. Image pixelization with differential privacy. In *Data and Applications Security and Privacy XXXII - 32nd Annual IFIP WG 11.3 Conference, DBSec 2018, Bergamo, Italy, July 16-18, 2018, Proceedings*, pages 148–162, 2018.
- [8] S. Fei, Z. Yan, W. Ding, and H. Xie. Security vulnerabilities of sgx and countermeasures: A survey. *ACM Comput. Surv.*, 54(6), 2021.
- [9] M. Fredrikson, S. Jha, and T. Ristenpart. Model inversion attacks that exploit confidence information and basic countermeasures. In *ACM Conference on Computer and Communications Security*, 2015.
- [10] M. Fredrikson, E. Lantz, S. Jha, S. Lin, D. Page, and T. Ristenpart. Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing. In *23rd USENIX Security Symposium USENIX Security*, pages 17–32, San Diego, CA, 2014. USENIX Association.
- [11] J. Gallier. *Geometric Methods and Applications for Computer Science and Engineering*. Springer-Verlag, New York, 2000.
- [12] 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 M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 201–210, 2016.
- [13] R. M. Heiberger. Generation of random orthogonal matrix. *Journal of the Royal Statistical Society*, 27(2), 1978.
- [14] B. Hitaj, G. Ateniese, and F. Perez-Cruz. Deep models under the gan: Information leakage from collaborative deep learning. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS '17*, pages 603–618, New York, NY, USA, 2017. ACM.
- [15] H. Hu, Z. Salicic, L. Sun, G. Dobbie, P. S. Yu, and X. Zhang. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37, 2022.
- [16] Y. Huang, Z. Song, K. Li, and S. Arora. InstaHide: Instance-hiding schemes for private distributed learning. In H. D. III and A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4507–4518. PMLR, 13–18 Jul 2020.
- [17] Z. Huang, W. jie Lu, C. Hong, and J. Ding. Cheetah: Lean and fast secure Two-Party deep neural network inference. In *31st USENIX Security Symposium (USENIX Security 22)*, pages 809–826, Boston, MA, Aug. 2022. USENIX Association.
- [18] M. Jagielski, N. Carlini, D. Berthelot, A. Kurakin, and N. Papernot. High accuracy and high fidelity extraction of neural networks. In *29th USENIX security symposium (USENIX Security 20)*, pages 1345–1362, 2020.
- [19] P. Kairouz and Others. Advances and open problems in federated learning. *Foundations and Trends in Machine Learning*, 14(1-2), 2021.
- [20] J. Katz and Y. Lindell. *Introduction to Modern Cryptography*. Chapman and Hall/CRC, 2007.
- [21] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Commun. ACM*, 60(6):84–90, May 2017.
- [22] G. B. H. E. Learned-Miller. Labeled faces in the wild: Updates and new reporting procedures. Technical Report UM-CS-2014-003, University of Massachusetts, Amherst, May 2014.
- [23] S. Lee, K. Chwa, J. Hahn, and S. Shin. Image morphing using deformation techniques. *JOURNAL OF VISUALIZATION AND COMPUTER ANIMATION*, 7(1):3 – 23, n.d.
- [24] M. Li, L. Lai, N. Suda, V. Chandra, and D. Z. Pan. Privynet: A flexible framework for privacy-preserving deep neural network training with A fine-grained privacy control. *CoRR*, abs/1709.06161, 2017.
- [25] S. Mansfield-Devine. The Ashley Madison affair. *Network Security*, 2015(9):8 – 16, 2015.
- [26] P. Mohassel and Y. Zhang. SecureML: A system for scalable privacy-preserving machine learning. In *2017 IEEE Symposium on Security and Privacy (SP)*, pages 19–38, 2017.
- [27] L. K. Ng, S. S. M. Chow, A. P. Y. Woo, D. P. H. Wong, and Y. Zhao. Goten: Gpu-outsourcing trusted execution of neural network training. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(17):14876–14883, May 2021.
- [28] E. Raff, J. Sylvester, S. Forsyth, and M. McLean. Barrage of random transforms for adversarially robust defense. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6521–6530, 2019.
- [29] D. Rathee, M. Rathee, N. Kumar, N. Chandran, D. Gupta, A. Rastogi, and R. Sharma. Cryptflow2: Practical 2-party secure inference. In *27th*

Annual Conference on Computer and Communications Security (ACM CCS 2020). ACM, 2020.

- [30] S. Sharma and K. Chen. Confidential boosting with random linear classifiers for outsourced user-generated data. In *Computer Security - ESORICS 2019 - 24th European Symposium on Research in Computer Security, Luxembourg, September 23-27, 2019, Proceedings, Part I*, pages 41–65, 2019.
- [31] S. Sharma, K. Chen, and A. Sheth. Toward practical privacy-preserving analytics for iot and cloud-based healthcare systems. *IEEE Internet Computing*, 22(2):42–51, Mar./Apr. 2018.
- [32] R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015.
- [33] R. Shokri and V. Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015.
- [34] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. Membership inference attacks against machine learning models. In *2017 IEEE Symposium on Security and Privacy, SP 2017, San Jose, CA, USA, May 22-26, 2017*, pages 3–18, 2017.
- [35] F. Tramer and D. Boneh. Slalom: Fast, verifiable and private execution of neural networks in trusted hardware. In *International Conference on Learning Representations*, 2019.
- [36] F. Tramèr, F. Zhang, A. Juels, M. K. Reiter, and T. Ristenpart. Stealing machine learning models via prediction apis. In *Proceedings of the 25th USENIX Conference on Security Symposium, SEC'16*, pages 601–618, USA, 2016. USENIX Association.
- [37] L. Unger. Breaches to customer account data. *Computer and Internet Lawyer*, 32(2):14 – 20, 2015.
- [38] L. van der Maaten and G. Hinton. Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9:2579–2605, 2008.
- [39] S. S. Vempala. *The Random Projection Method*. American Mathematical Society, 2005.
- [40] P. Xie, M. Bilenko, T. Finley, R. Gilad-Bachrach, K. E. Lauter, and M. Naehrig. Crypto-nets: Neural networks over encrypted data. *CoRR*, abs/1412.6181, 2014.
- [41] Q.-s. Zhang and S.-c. Zhu. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology and Electronic Engineering*, 2018.
- [42] Y. Zhang, R. Jia, H. Pei, W. Wang, B. Li, and D. Song. The secret revealer: Generative model-inversion attacks against deep neural networks. In *CVPR*, 2020.