Toward Secure and Efficient Deep Learning Inference in Dependable IoT Systems

Han Qiu[®], Member, IEEE, Qinkai Zheng, Tianwei Zhang[®], Meikang Qiu[®], Senior Member, IEEE, Gerard Memmi, Member, IEEE, and Jialiang Lu

Abstract—The rapid development of deep learning (DL) enables resource-constrained systems and devices [e.g., Internet of Things (IoT)] to perform sophisticated artificial intelligence (AI) applications. However, AI models, such as deep neural networks (DNNs), are known to be vulnerable to adversarial examples (AEs). Past works on defending against AEs require heavy computations in the model training or inference processes, making them impractical to be applied in IoT systems. In this article, we propose a novel method, SUPER-IOT, to enhance the security and efficiency of AI applications in distributed IoT systems. Specifically, SUPER-IOT utilizes a pixel drop operation to eliminate adversarial perturbations from the input and reduce network transmission throughput. Then, it adopts a sparse signal recovery method to reconstruct the dropped pixels and wavelet-based denoising method to reduce the artificial noise. SUPER-IOT is a lightweight method with negligible computation cost to IoT devices and little impact on the DNN model performance. Extensive evaluations show that it can outperform three existing AE defensive solutions against most of the AE attacks with better transmission efficiency.

Index Terms—Adversarial examples (AEs), deep learning (DL), Internet of Things (IoT), security.

I. INTRODUCTION

THE PAST decade has witnessed the revolutionary development of deep learning (DL) technology with deep neural networks (DNNs). A variety of DL algorithms and models were designed to perform different artificial intelligence (AI) tasks. For instance, convolutional neural networks (CNNs) [1] show great capability in handling computer vision tasks; recurrent neural networks (RNNs) [2] power the advance of natural language processing; and deep reinforcement learning (DRL) [3] achieves very high performance in robotics and

Manuscript received April 30, 2020; revised June 5, 2020; accepted June 19, 2020. Date of publication June 23, 2020; date of current version February 19, 2021. This work was supported by the Open Research Project of the State Key Laboratory of Industrial Control Technology, Zhejiang University, China under Grant ICT 20025. (*Corresponding author: Meikang Qiu.*)

Han Qiu, Qinkai Zheng, and Gerard Memmi are with Telecom Paris, Institut Polytechnique de Paris, 91120 Palaiseau, France (e-mail: han.qiu @telecom-paris.fr; qinkai.zheng@telecom-paris.fr; gerard.memmi @telecom-paris.fr).

Tianwei Zhang is with the School of Computer Science and Engineering, Nanyang Technological University, 639798 Singapore (e-mail: tianwei.zhang@ntu.edu.sg).

Meikang Qiu is with the Department of Computer Science, Texas A&M University–Commerce, Commerce TX 75428 USA (e-mail: qiumeikang@yahoo.com).

Jialiang Lu is with SPEIT, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: jialiang.lu@sjtu.edu.cn).

Digital Object Identifier 10.1109/JIOT.2020.3004498

autonomous driving. Those state-of-the-art models have been extensively commercialized in many products, and they are continuously enhanced by experts from academia as well as industry. Nowadays, new techniques have kept emerging at surprising speed to enrich the DL community.

Meanwhile, DL also drives the growth of the Internet of Things (IoT). Equipped with different sensors (e.g., cameras, microphones, and gyroscopes), IoT devices become appealing targets for DL applications. They keep sensing data and information from various environmental contexts in a streaming fashion. Then, DL models are deployed in centralized servers to process and understand the data. The integration of AI and IoT leads to the era of AI of Things (AIoT), which have significantly changed our daily life (Fig. 1): small-scaled AIoT systems are introduced to build smart homes and increase the comfort and quality of life; medium-scale AIoT systems are deployed in warehouses and factories for higher efficiency and automation; and large-scale AIoT systems can contribute to the establishment of smart cities.

Two challenges need to be addressed for the deployment of DL models in the AIoT systems. The first one is *efficiency*. An IoT system can consist of a large number of edge devices with high-quality sensors streaming information at a very high rate (e.g., remote sensing [4]). This can result in a large amount of data transferring between the sensor devices and the model server [5]. There could be a performance bottleneck if the network bandwidth of the AIoT system is not high, or an energy bottleneck if the transmission energy budget is low. Thus, it is necessary to have an efficient approach to processing the data at the sensor device before sending them to the model server, in order to reduce the throughput and transmission energy dissipation.

The second challenge we need to consider is *security*. DL models are well known to be vulnerable to adversarial examples (AEs) [6]. An AE is created by adding imperceptible perturbations to a clean data sample, which can mislead the model to give a wrong decision. Past works have demonstrated that an adversary can generate such AEs of images [7], voices [8], and laser signals [9] to spoof the IoT sensors and cause catastrophic consequences. It is of paramount importance for the sensor devices to detect or prevent such malicious samples for secure model inference.

To the best of our knowledge, currently, there are no existing solutions that can solve both of the two challenges. The most promising direction is to add a preprocessing step on the input samples before feeding them into the model [10]–[12].

2327-4662 © 2020 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.



Fig. 1. Different scales of AIoT systems.

Such a step introduces nondifferentiable transformations on the inputs to obfuscate the gradients of the models, so the difficulty of AE generation is increased and the impact of the calculated perturbations is mitigated. However, such defense approaches are still vulnerable as the adversary can adaptively and statistically calculate the gradient based on the preprocessing algorithms [13]. Besides, some preprocessing operations can introduce heavy computation (e.g., sophisticated quantization in image compression [12]), which are not applicable to computing resource-constrained IoT devices.

In this article, we propose SUPER-IOT: a secure and efficient approach to DL inference for dependable IoT systems, to overcome the two challenges.

- The essential component of our methodology is a pixel drop operation on the IoT ends, which randomly selects and drops a certain amount of pixels of the input images. Such operation can reduce the data throughput between the IoT device and server to achieve higher network efficiency. At the same time, it also gets a high chance to invalidate the effects of AEs since it could drop the added perturbations.
- 2) It is worth noted that the pixel drop operation can affect the model accuracy, especially for the clean samples, as it removes certain information which can be critical for model prediction. To maintain high performance, we adopt a novel pixel reconstruction algorithm, sparse signals recovery, on the model server to recover the dropped pixels.
- 3) To further enhance the performance of the model on adversarial as well as clean samples, we integrate a wavelet-based denoising operation on the model server to remove the adversarial perturbations and artificial noises.

We conducted extensive evaluations to demonstrate the effectiveness of SUPER-IOT. For security, we measured the defense effects of our solution against six state-of-the-art adversarial attacks. We also compared SUPER-IOT with three existing defense methods (Shield [12], pixel deflection (PD) [10], and feature distillation (FD) [11]): SUPER-IOT can maintain higher model accuracy and lower attack success rate than most defenses. For network efficiency, we measured the size of bitstreams with our preprocessing approach. SUPER-IOT can effectively reduce as high as 25% transmission throughput, while past works can hardly optimize network efficiency.

The major contributions of this article include: 1) a pixel drop operation to reduce the network throughput and mitigate AEs; 2) a pixel reconstruction algorithm to recover the original input and maintain high model accuracy; and 3) a wavelet-based denoising operation to further remove the adversarial perturbations and artificial noises.

This article is organized as follows. Section II discusses the research background and related works. Section III presents the problem definition and threat model. Section IV describes the design details of SUPER-IOT. Section V presents the evaluation results. We conclude in Section VI.

II. BACKGROUND AND RELATED WORKS

In this section, we briefly present the background and relevant works about AIoT systems, adversarial attacks on DNN models, and preprocessing-based defensive strategies.

A. Artificial Intelligence of Things

Benefiting from the advance of the DL technology, IoT systems are becoming more intelligent and multifunctional. A typical IoT network can consist of an enormous amount of IoT devices. They are connected via different communication technologies (e.g., Ethernet and Wi-Fi). They collect sensory data over time and transmits them to one or more centralized hosts. These hosts can be remote cloud servers, local gateways, or powerful edge devices. They run the DNN inference applications, interpret the received sensory data, and make control decisions. Such IoT configuration has been widely adopted in many scenarios, such as face authentication [14], vehicle detection [15], and remote monitoring [4].

The sensory data generated from the IoT devices exhibit some unique features. First, there can be a large quantity of connected IoT devices generating real-time data continuously. This leads to a huge volume of streaming data in the network. Second, various IoT devices can collect different types of sensory data and information, resulting in data heterogeneity. Those data need to be transmitted to the DNN inference engine and processed promptly to extract immediate insights and make fast decisions. These requirements need to be achieved from different perspectives.

- 1) At the host level, we can utilize powerful cloud servers with high computing capability and execution parallelism, or specialized hardware circuits [16] to accelerate the DNN inference.
- At the DNN algorithm level, novel algorithms were proposed (e.g., OS-ELM [17] and Faster R-CNN [18]) to handle the data streaming for object detection and video analytics.
- 3) At the network level, one possible method is to preprocess and compress the sensory data to reduce network throughput and communication costs. This is also what we aim to optimize.

B. Adversarial Attacks on DNN models

An adversary can add human-unnoticeable perturbations on the original input to fool a DNN classifier. Formally, as in (1), the target DNN model is a mapping function *F*. Given a clean input sample *x*, the corresponding AE is denoted as $\tilde{x} = x + \delta$, where δ is the adversarial perturbation. δ is constrained by certain metric (e.g., L_p norm) to make it imperceptible. Then, AE generation can be formulated as the optimization problem in (1a) (targeted attack where $l' \neq F(x)$ is the desired label set by the attacker, e.g., a cat image is misclassified specifically as a dog) or (1b) (untargeted attack, e.g., a cat image is misclassified as an arbitrary class other than a cat). In this article, we only evaluate the defense against the targeted attack and the untargeted attack can be mitigated in the same way

$$\min \|\|\delta\|, \text{ s.t. } F(\widetilde{x}) = l' \tag{1a}$$

min
$$\|\delta\|$$
, s.t. $F(\tilde{x}) \neq F(x)$. (1b)

Various approaches were proposed to solve the optimization problem and generate AEs. The fast gradient sign method (FGSM) [6] calculates the sign of the gradient of the classification loss with respect to the input sample, which gives the direction to modify input pixel values under L_{inf} constraints to generate AEs. Later on, variations of FGSM were introduced to iteratively calculate the perturbations with a small step or with momentum, e.g., I-FGSM [7]. Some approaches use a more advanced optimization algorithm to find the minimal adversarial perturbation under the L_2 constraint, such as LBFGS [19], DeepFool [20], and Carlini and Wagner (CW) [21].

Attack Scenarios: Generally, there are three attack scenarios [22], determined by the adversary's knowledge level of the target DNN model.

- 1) *White-Box Scenario:* The adversary knows every detail about the model including all the parameters. He can directly adopt the above approaches to generate AEs.
- Black-Box Scenario: The adversary does not have any knowledge about the target model. He has to use an alternative model of the same task to generate AEs and attack the target one.
- 3) Gray-Box Scenario: The adversary knows all details of the model (e.g., training algorithms, network topology, and hyperparameters) except the parameters. He can train another similar model with the same configurations for AE generation. The transferability property of AEs [23] can guarantee high success rates for black-box and gray-box scenarios.

C. Preprocessing-Based Defenses Against AEs

Various defensive strategies have been designed to defeat adversarial attacks. One direction is to train a more robust model from either scratch or an existing model. Those approaches aim to rectify AEs' malicious features by including AEs into the training set [24], processing all the training data [25], or revising the DNN topology [26]. However, training a DNN model is very time and resource consuming, especially when the model is complicated. Besides, those methods are not applicable when the DNN models are packed as closed-source applications and cannot be modified. Most of all, those methods are not secure: the adversary can still generate adaptive AEs for the new models [27].

A more promising direction is to preprocess the input data to eliminate adversarial influence without touching the DNN model. Typical transformation methods include denoising, compression, drop pixels, etc. These solutions are suitable in AIoT systems, as it is feasible and efficient to preprocess the sensory data on IoT devices. So here, we focus on this preprocessing-based direction. Below, we describe some existing works and their limitations. We introduce our novel preprocessing technique in Section IV, and empirically demonstrate its advantages over those works in Section V.

Shield [12]: In this approach, JPEG compression is improved by randomizing the quantization factors to different blocks of image contents. Then, the compression process consists of the discrete cosine transform (DCT) and lossy quantization. This nondifferentiable and irreversible transformation can obfuscate the gradients of the DNN model with respect to inputs from the adversary. However, this method can also decrease the classification accuracy of clean samples.

Feature Distillation [11]: This approach uses a revised JPEG compression-based mechanism to defeat AEs. The quantization step in the DCT process is modified to optimize the reduction of the adversarial perturbations to improve the robustness of the DNN model. However, FD is inefficient as this revised quantization step can reduce the compression ratio.

Pixel Deflection [10]: The idea of this approach is to combine the denoising algorithm with the operation of dropping pixels. First, around 0.1% pixels of the input image is dropped and replaced with a random pixel value within a small range. Then, the denoising technique is applied to reduce the adversarial perturbations. PD could provide robustness against the AEs. But the compression ratio (i.e., dropped pixels) has to be very small in order to maintain the model's prediction accuracy.

III. PROBLEM DEFINITION AND THREAT MODEL

In this article, we aim to design a novel methodology for secure and efficient DNN inference for AIoT systems. Specifically, we consider a distributed IoT system conducting computer vision tasks (e.g., image classification and object detection). The sensor devices in the system keep collecting the visual input at high sampling rates and sending them to a centralized server for DNN inference. We focus on computer vision applications for two reasons. First, vision sensors (e.g., cameras) are one of the most widely used IoT devices in our daily life. Computer vision tasks are also commonly adopted in many scenarios, e.g., video surveillance [15], face authentication [14], autonomous driving [28], etc. Second, compared with other sensory data and tasks, vision sensors can produce a larger volume of real-time streaming data with higher throughput. So it is in a more urgent need for an efficient inference solution in an IoT system.

We assume that the DNN model deployed in the IoT system cannot be modified. In reality, the IoT administrator can purchase the DNN model from a model vendor. He may not be allowed to customize the model due to intellectual property protection. He may not be able to alter the model either if it is packed as a closed-source application. Then past approaches to training or retraining models for better robustness cannot be applied in our case. Designing new DNN hardware accelerators for better performance is out of the scope of this article, as it requires drastic changes to the underlying infrastructure with high cost. The administrator can only implement preprocessing



Fig. 2. Methodology overview.

functions on the IoT devices or the server. Those functions must meet the following requirements.

- 1) *Efficiency:* They must be able to reduce the throughput of the transmitted data from the IoT devices to the DNN server to relieve the burden and stress of the network bandwidth.
- Lightweight: The preprocessing function on the IoT devices should not be too heavy to impact the devices' performance or operations, considering the limited onboard computing capabilities and resources.
- Functionality Preserving: They should not affect the prediction accuracy of the DNN model on clean data samples.
- 4) *Security:* They should be able to effectively remove the adversarial perturbations and preserve the correct prediction results from the model.

Threat Model: We assume the entire IoT system (e.g., sensor devices, central servers, and communication channel) is trusted. So we do not consider the security threats from IoT botnet (e.g., Mirai [29] and Hajime [30]) and man-in-themiddle attacks. We also assume the target DNN model is correct without DNN backdoors [31]. The adversary is outside of the IoT system, attempting to spoof the sensors and DNN model by adding malicious perturbations on the physical objects or spots on the lens of the cameras [32]. He has white-box access to the DNN model and the preprocessing functions. We aim to show that even the adversary knows all detail of the target model and the defense mechanism, he can still not generate AEs to bypass the defense to compromise the model.

IV. PROPOSED METHODOLOGY

In this section, we present our efficient and secure approach for DNN inference in IoT systems. We give the methodology overview in Section IV-A, following by the descriptions of each operation in Sections IV-B–IV-D.

A. Design Overview

Fig. 2 shows an overview of our proposed methodology. It consists of three steps across the IoT device and model server. The first step is conducted on the IoT device, which randomly drops some pixels from the input image. This operation can remove the potential adversarial perturbations with a high chance if the drop rate is high. Meanwhile, it can also increase the transmission efficiency between the IoT device and the model server as the image size is reduced after dropping certain pixels. This operation is lightweight and incurs very little computing costs on the IoT device. After this operation, the image will be sent out to the model server.

On the server, the received image cannot be directly fed into the DNN model, as a lot of information has been removed. Then, the second step is to reconstruct the dropped pixels. This operation can approximately recover the dropped pixels other than the perturbations. It will increase the model's prediction accuracy on this image.

The last step is image denoising. This operation can remove the malicious perturbations that are not dropped out at the IoT side, and also the artificial noises introduced during the reconstruction process. After that, the image can be sent to the DNN model for classification. Below, we detail the mechanism and algorithm of each step.

B. Step 1: Pixel Dropping

In this step, the IoT device randomly selects a fixed ratio r of pixels and remove them out of the image. This can reduce the transmission throughput, and also remove adversarial perturbations. Specifically, we first divide the raw image into multiple blocks of $N \times N$ pixels. We denote one block as $f_0(x, y)$, where (x, y) represents the coordinate of pixels. Then, we randomly select $n = r \times N \times N$ pixel inside each block and set their values as zero. As shown in (2), n pixels $\{(x_1, y_1), \ldots, (x_n, y_n)\}$ are set to zero, while the rest remains the same as $f_0(x, y)$. The resulting block is denoted as $f_1(x, y)$. Finally, we concatenate the new blocks into one output image and send it to the model server

$$f_1(x, y) = \begin{cases} 0, & \text{if } (x, y) \in \{(x_1, y_1), \dots, (x_n, y_n)\} \\ f_0(x, y), & \text{otherwise.} \end{cases}$$
(2)

Note that the value of r can determine the efficiency, security, and also the model performance: a large r can reduce more throughput and decrease the success rate of adversarial perturbations. However, it can also decrease the model performance on clean samples. So we must carefully select r to balance such a tradeoff. Fig. 3 (first row) shows the output images with different drop ratios. We will empirically identify the optimal value in Section V.

C. Step 2: Pixel Reconstruction

When receiving the compressed image, the model server adopts the pixel reconstruction algorithm inspired by sparse signals recovery [33]. The algorithm is processed block by block using 2D-DCT transform

$$F(u, v) = \frac{2}{N} \sum_{i=1}^{N} \sum_{j=0}^{N} \alpha_{i,j}(u, v) f(i, j)$$

$$\alpha_{i,j}(u, v) = \Lambda(i)\Lambda(j) \cos\left[\frac{\pi u}{2N}(2i-1)\right] \cos\left[\frac{\pi v}{2N}(2j-1)\right]$$

$$\Lambda(x) = \begin{cases} \frac{1}{\sqrt{2}}, & \text{if } x = 0\\ 1, & \text{otherwise.} \end{cases}$$
(3)



Fig. 3. Visual content evaluation of pixel drop with different ratio (r from 0.1 to 0.7) and the reconstruction results measured by PSNR.



Fig. 4. Dynamic process of image reconstruction. (a) (MSE and PSNR) versus iterations. (b) $[\Delta, \mu]$ versus iterations.

The reconstruction algorithm is detailed in Algorithm 1: given a block f_1 , for each dropped pixel at position (k, l) $((k, l) \in \{(x, y) | f_1(x, y) = 0\})$, we estimate a gradient to modify its pixel value. We first perturb the pixel value in two directions with a distortion level Δ (lines 4 and 5). Then, we calculate their 2D-DCT transform and L_1 norm, respectively (lines 6-9). The gradient is calculated (line 10) and used to update the pixel in block f_1 with a step size of μ (line 11). During this iterative process, we keep monitoring the changes of reconstructed images using the metric mean-square error (MSE), which is defined as the difference of output images in two consecutive iterations (line 13). When MSE is hardly changed, we dynamically reduce Δ and μ to achieve better reconstruction results (lines 14-18). This iterative process will end until MSE is smaller than a threshold ϵ (line 21).

Fig. 4 shows the trends of MSE, peak signal-to-noise ratio (PSNR), and $[\Delta, \mu]$ during the image reconstruction process. Initially, MSE keeps decreasing while PSNR keeps increasing. When MSE becomes too small, the reconstruction tends to converge and stops the modification of pixel values. To refine the reconstruction, when MSE is smaller than 1% of the maximal MSE previously, Δ and μ are updated dynamically (on the 13th and 17th iterations). At the 20th iteration, MSE already becomes very small and the reconstructed image has a good quality with PSNR bigger than 30. Then, the iterative process stops.

The hyperparameters used in this pixel reconstruction algorithm can significantly affect the difficulty of the reconstruction process and the quality of the output. So we need to discover the optimal values.

First, for the initial values of distortion level Δ and step size μ , we tried different values within the range [0.01, 0.1]

Algorithm 1: Pixel Reconstruction

Input: a image block $f_1(x, y)$ **Output**: the reconstructed image block $f_2(x, y)$ **Parameters**: Δ distortion level; μ step size; ϵ stop criterion.

/* Initialization */ -0

$$d = 0, f^{0}(x, y) = f_{1}(x, y), MSE_{max} = 0;$$

1

7

8

9

11

12

13

17

19

20

3 **for** (k, l) in $\{(x, y)|f_1(x, y) = 0\}$ **do** /* Perturb pixel value of missing pixel; δ Dirac function */ $f_{+}^{(k,l)}(x,y) = f_{1}(x,y) + \Delta\delta(x-k,y-l);$ 4 $f_{-}^{(k,l)}(x, y) = f_1(x, y) - \Delta \delta(x - k, y - l);$ /* 2D-DCT transform */ 5
$$\begin{split} F_{+}^{(k,l)}(u,v) &= \frac{2}{N} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_{i,j}(u,v) f_{+}^{(k,l)}(i,j); \\ F_{-}^{(k,l)}(u,v) &= \frac{2}{N} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_{i,j}(u,v) f_{-}^{(k,l)}(i,j); \end{split}$$
6 /* L1 norm */
$$\begin{split} \|F_{+}^{(k,l)}\|_{1} &= \sum_{u=1}^{N} \sum_{v=1}^{N} \|F_{+}^{(k,l)}(u,v)\|_{1}; \\ \|F_{-}^{(k,l)}\|_{1} &= \sum_{u=1}^{N} \sum_{v=1}^{N} \|F_{-}^{(k,l)}(u,v)\|_{1}; \end{split}$$
/* Estimate gradient */ $grad(k, l) = \frac{\|F_{+}^{(k,l)}\|_{1} - \|F_{-}^{(k,l)}\|_{1}}{2\Delta};$ /* Update pixel value */ 10 $f^{d+1}(k, l) = f^d(k, l) - \mu \times grad(k, l);$ end /* Dynamically update Δ , μ */ $MSE = \|f^{d+1}(x, y) - f^d(x, y)\|_2;$ if $MSE < 0.01 \times MSE_{max}$ then 14 $\Delta = \Delta/10;$ 15 $\mu = \mu / 10;$ 16 $MSE_{max} = 0;$ 18 end $MSE_{max} = max(MSE, MSE_{max});$ d = d + 1;21 while $MSE > \epsilon$; 22 $f_2(x, y) = f^d(x, y);$ 23 return $f_2(x, y)$

and measured the quality of reconstructed images using the metric PSNR, which is defined as the visual content deviation from the clean image. Fig. 5 shows the average PSNR for each



Fig. 5. PSNR of reconstructed images under different values of Δ and μ .

configuration. We can choose $\Delta = 0.03$ and $\mu = 0.02$ that lead to the best image quality.

Second, the stop criterion ϵ determines the number of iterations during image reconstruction. As shown in Fig. 4, the MSE and PSNR will become saturated after a certain number of rounds. Then, it is not necessary to continue the iteration, as the quality of the reconstructed image will not change. So an appropriate threshold ϵ can guarantee the best quality of output with the minimal number of iterations. We empirically identify the optimal $\epsilon = 10^{-5}$ from Fig. 4.

Third, the pixel drop ratio r can also impact the image reconstruction. Fig. 3 (second row) shows the reconstructed images and their PSNRs with different drop ratios. We observe that larger r leads to smaller PSNR (i.e., worse quality). Besides, the value of r can also determine the effects of AEs and network throughput. More evaluation results will be presented in Section V to show the tradeoff between those aspects, and discover the ideal drop ratio.

D. Step 3: Image Denoising

After the pixel reconstruction, the model server uses the image denoising algorithm to further improve the image quality. On the one hand, the pixel dropping and reconstruction can introduce artificial noises. Then this denoising operation can filter such new introduced noises [10]. On the other hand, this denoising operation is nondifferentiable. It can obfuscate the DNN model gradients to further increase the difficulty of AE generations via gradient-based approaches.

In this article, we adopted the wavelet-based denoising method named BayesShrink [34] from [10]. Other denoising methods can be applied in a similar way. The denoising method is performed in the frequency domain through the wavelet transform. The image noise is always assumed as the small perturbations on values in the high-frequency domain. Therefore, these small values can be removed by setting coefficients below a given threshold to zero (hard threshold) or shrinking different coefficients toward zero by a soft threshold. First, we use the VisuShrink approach to set a hard threshold. For an image X with N pixels, this threshold is

given by $\sigma \sqrt{2\log N}$, where σ is normally smaller than the true noise standard deviation. Then, we adopted the BayesShrink algorithm [10] as an additional step to set a soft threshold to further filter the wavelet coefficients. The threshold $T_h * (\sigma_x, \beta)$ is estimated on each wavelet sub-band and the optimal threshold is calculated by minimizing the expected MSE. We model the threshold for each wavelet coefficient as a generalized Gaussian distribution (GGD). It can be approximated as (σ^2/σ_x) , where σ_x and β are parameters of the GGD for each wavelet sub-band

$$T_h * (\sigma_x, \beta) = \underset{T_h}{\operatorname{argmin}} E(\widehat{X} - X)^2 \approx \frac{\sigma^2}{\sigma_x}.$$
 (4)

Normally, an approximation of T_h , as shown on the righthand side of (4), is used to adapt to the amount of noise in the given image. The parameters for the denoising in this article are tuned to get the best performance.

V. EVALUATIONS

In this section, we comprehensively evaluate the efficiency and security of our proposed methodology. We measure its resilience against six popular adversarial attacks and compare it with three existing preprocessing-based defense methods. We also measure and compare the network throughput benefits introduced by different approaches.

A. Experimental Configuration

We consider an image classification task on the CIFAR-10 data set. There are 50 000 images for training and 10 000 images for testing. Each image has a size of $32 \times 32 \times 3$ and belongs to one of ten classes. All pixel values are normalized within the range of [0, 1].

We choose ResNet-29 [35] as the target model. It consists of 29 layers for three bottleneck residual blocks with channel sizes of 64, 128, and 256, respectively. We use the Keras package with Tensorflow 1.14 [36] backend to implement the model. Weights in all convolutional layers are initialized by a truncated normal distribution proposed in [37]. The training process is done via the Adam optimization algorithm [38] with its hyperparameters $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The model is trained to reach the top-1 accuracy of 92.27% over the testing set after about 150 epochs. Experiments are conducted on a server with a CPU of Intel Core i9-9900K@3.60 GHz and a GPU of NVIDIA GeForce GTX 2080 Ti.

Hyperparameters: For the pixel drop and reconstruction algorithms, we set N = 8 to have blocks with 8×8 pixels, which is a typical configuration of DCT transform applied in image compression. We set $\Delta = 0.03$, $\mu = 0.02$, and $\epsilon = 10^{-5}$, as discussed in Section IV-C.

B. Security Evaluation

First, we check whether our methodology can defeat existing adversarial attacks. We consider six well-known attack techniques: 1) FGSM [19]; 2) I-FGSM [7]; 3) DeepFool [20]; 4) LBFGS [6]; 5) CW [21]; and 6) PGD [39]. We adopt the CleverHans library (v3.0.1) [40] to generate AEs with those approaches. We set the $|L_2|$ between AEs and original images

 TABLE I

 TOP-1 ACCURACY IN THE PRESENCE OF VARIOUS ADVERSARIAL ATTACKS ON SUPER-IOT WITH DIFFERENT PIXEL DROP RATIOS

Attack	Baseline	r=0.01	r=0.05	r=0.10	r=0.15	r=0.20	r=0.25	r=0.30	r=0.40	r=0.50
Clean	1.00	1.00	0.99	0.99	0.93	0.94	0.88	0.86	0.60	0.43
FGSM	0.39	0.60	0.65	0.70	0.72	0.75	0.81	0.79	0.59	0.41
I-FGSM	0.21	0.45	0.54	0.65	0.78	0.79	0.79	0.71	0.64	0.42
PGD	0.05	0.30	0.39	0.51	0.67	0.69	0.72	0.72	0.56	0.51
DeepFool	0.00	0.99	0.98	0.97	0.92	0.89	0.85	0.83	0.60	0.44
LBFGS	0.00	0.95	0.99	0.97	0.92	0.89	0.85	0.82	0.59	0.37
CW	0.00	0.98	0.97	0.96	0.93	0.89	0.88	0.80	0.62	0.43

to be within 0.5, to make the perturbation imperceptible. For FGSM and I-FGSM, the scale of distortion is $\epsilon = 0.005$ under the Linf constraints. For PGD, the scale of distortion is $\epsilon = 0.01$ and the number of attack iterations is 10. For the rest of the attacks, the optimization process is iterated until the adversary generates AEs of all samples. For all evaluations, we consider the targeted attack, where a random label different from the correct one is selected as the adversary's target. The AEs are generated under a white-box scenario. Table I shows the prediction accuracy of the clean samples as well as AEs when they are not preprocessed (Baseline column), or preprocessed by our methodology with different drop ratio r. All average accuracy is measured for 100 images. We observe that without any defense, all attacks can significantly compromise the performance of the target model, even making the accuracy drop to 0. With our preprocessing operation, the prediction accuracy is significantly increased. The accuracy of classifying AEs from FGSM, I-FGSM, and PGD will increase first and then decrease which is different from the DeepFool, LBFGS, and CW. This is due to the initial value r = 0.01 has been already effective to mitigate the DeepFool, LBFGS, and CW attacks. Even lower r will see the same accuracy trend of DeepFool, LBFGS, and CW compared with FGSM, I-FGSM, and PGD. Then, continuing increase r after 0.25 will lead to the decrease of model accuracy since and is a tradeoff between classifying clean samples and mitigating AEs.

Next, we compare SUPER-IOT with existing state-of-theart solutions: Shield [12], PD [10], and FD [11]. Those preprocessing-based solutions do not require the modification of DNN models and can be applied to our IoT scenario. Due to the stochastic features in these solutions, we repeat the experiments ten times and report the average prediction accuracy of each preprocessing method for each attack technique, shown in Table II. (We set r = 0.1 in our method.) We observe that our solution can beat the other methods on the performance of AEs from DeepFool, LBFGS, and CW. For AEs from FGSM, I-FGSM, and PGD, the accuracy is higher than PD but slightly lower than Shield and FD. However, Shield and FD have bad performance on the clean samples, making them less practical.

C. Efficiency Analysis

We measure the efficiency of SUPER-IOT in terms of reduced transmission throughput. By dropping certain pixels, SUPER-IOT can effectively reduce the total bitstream for transmission. This can save the energy cost of the IoT devices, and relieve the stress of network bandwidth in IoT systems.

TABLE II TOP-1 ACCURACY IN THE PRESENCE OF VARIOUS ADVERSARIAL ATTACKS ON THE BASELINE MODEL, SHIELD, FD, PD, AND SUPER-IOT

Attack	$L_{ m inf}$	L_2	Baseline	Shield	FD	PD	SUPER
Clean	0.0	0.0	1.00	0.85	0.94	0.98	0.99
FGSM	0.005	0.28	0.39	0.78	0.85	0.57	0.70
I-FGSM	0.005	0.21	0.21	0.79	0.83	0.44	0.65
PGD	0.010	0.39	0.05	0.79	0.70	0.29	0.51
DeepFool	0.015	0.12	0.00	0.82	0.91	0.87	0.97
LBFGS	0.017	0.15	0.00	0.82	0.92	0.97	0.97
CW	0.115	0.09	0.00	0.82	0.93	0.94	0.96



Fig. 6. Evaluation of SUPER-IOT on the DNN model classification accuracy (on both clean sample and AEs generated by two different approaches) and compression ratio with different pixel drop ratio.

Fig. 6 shows the compression ratio evaluated of SUPER-IOT with different pixel drop ratio ranging from 0.01 to 0.5 (the bars). The evaluation here is made based on the bitstream size before feeding into the compression algorithms. We can see a larger r leads to a larger compression ratio for better efficiency. This is straightforward: when r of the pixels are dropped out, the size of transmitted data will also be reduced by r. However, as we discussed in Section V-B, a larger rcan also affect the prediction accuracy of clean samples and AEs (curves in Fig. 6). So users need to carefully balance the tradeoff between efficiency, security, and functionality, and consider their requirements when configuring the drop ratio.

In contrast, past works on AE defenses cannot achieve throughput reduction. For PD, the optimal pixel drop ratio is between 0.1% and 1% in order to have a good model performance for both adversarial and clean images. This ratio has very little improvement in transmission efficiency. For Shield and FD, there are no data reduction effects at the bitstream level in their preprocessing operations. Thus, considering the transmission efficiency, SUPER-IOT has an obvious advantage over PD, Shield, and FD. For the computation overhead, both PD and our SUPER-IOT only require pixel drop operation on the IoT devices which will add very limited additional computing overhead. For comparison, Shield and FD have operations of modifying the JPEG compression process with sophisticated randomization or dedicated quantization tables, which are much heavier than the pixel drop operation for IoT devices.

We adopt state-of-the-art configuration, where IoT collects sensory data, and sends them to the server for inference. The only computation cost introduced is the image reconstruction at the server end, which is lightweight. However, we can save the network throughput and transmission power from the IoT end. Since IoT devices are more resource constrained than the server, such optimization is meaningful.

VI. FUTURE WORK AND CONCLUSION

For future work, we plan to explore better data drop and reconstruction methods for higher classification accuracy on both clean samples and AEs. With a larger drop ratio and better reconstruction methods, the network throughput can be further reduced. We will also explore more advanced defense solutions for AIoT systems against more advanced attacks like adaptive adversarial attacks.

In this article, we proposed a novel approach, SUPER-IOT, to efficiently secure the inference of DNN models in AIoT systems. We employ three techniques (pixel dropping, pixel reconstruction, and image denoising) to defeat AEs and maintain good performance for clean samples. Meanwhile, those operations can also achieve high efficiency for network transmission in the IoT systems by reducing the bitstream of transmitted sensory data. Our approach is lightweight with little impact on the IoT devices' performance or operations. It is generic and can be applied to various computer vision tasks without modifying the DNN models.

REFERENCES

- A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran Assoc., Inc, 2012, pp. 1097–1105.
- [2] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, Vancouver, BC, Canada, 2013, pp. 6645–6649.
- [3] V. Mnih *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [4] G. Chen, T. X. Han, Z. He, R. Kays, and T. Forrester, "Deep convolutional neural network based species recognition for wild animal monitoring," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Paris, France, 2014, pp. 858–862.
- [5] H. Qiu, Q. Zheng, G. Memmi, J. Lu, M. Qiu, and B. Thuraisingham, "Deep residual learning based enhanced JPEG compression in the Internet of Things," *IEEE Trans. Ind. Informat.*, early access, May 14, 2020, doi: 10.1109/TII.2020.2994743.
- [6] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," 2014. [Online]. Available: arXiv:1412.6572.
- [7] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial examples in the physical world," 2016. [Online]. Available: arXiv:1607.02533.
- [8] N. Carlini et al., "Hidden voice commands," in Proc. 25th USENIX Security Symp. (USENIX Security), 2016, pp. 513–530.
- [9] Y. Cao et al., "Adversarial sensor attack on lidar-based perception in autonomous driving," in Proc. ACM SIGSAC Conf. Comput. Commun. Security, 2019, pp. 2267–2281.

- [10] A. Prakash, N. Moran, S. Garber, A. DiLillo, and J. Storer, "Deflecting adversarial attacks with pixel deflection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Salt Lake City, UT, USA, 2018, pp. 8571–8580.
- [11] Z. Liu et al., "Feature distillation: DNN-oriented JPEG compression against adversarial examples," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Long Beach, CA, USA, 2019, pp. 860–868.
- [12] N. Das et al., "Shield: Fast, practical defense and vaccination for deep learning using JPEG compression," in Proc. 24th ACM SIGKDD Int. Conf. Knowl. Discover. Data Mini., 2018, pp. 196–204.
- [13] A. Athalye and N. Carlini, "On the robustness of the CVPR 2018 white-box adversarial example defenses," 2018. [Online]. Available: arXiv:1804.03286.
- [14] Y. Mao, S. Yi, Q. Li, J. Feng, F. Xu, and S. Zhong, "A privacy-preserving deep learning approach for face recognition with edge computing," in *Proc. USENIX Workshop Hot Topics Edge Comput. (HotEdge)*, 2018, pp. 1–6.
- [15] Y. Tang, C. Zhang, R. Gu, P. Li, and B. Yang, "Vehicle detection and recognition for intelligent traffic surveillance system," *Multimedia Tools Appl.*, vol. 76, no. 4, pp. 5817–5832, 2017.
- [16] S. Han *et al.*, "EIE: efficient inference engine on compressed deep neural network," *ACM SIGARCH Comput. Archit. News*, vol. 44, no. 3, pp. 243–254, 2016.
- [17] Z. Yang, P. Zhang, and L. Chen, "RFID-enabled indoor positioning method for a real-time manufacturing execution system using OS-ELM," *Neurocomputing*, vol. 174, pp. 121–133, Jan. 2016.
- [18] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards realtime object detection with region proposal networks," in *Advances in Neural Information Processing Systems*. Red Hook, NY, USA: Curran, 2015, pp. 91–99.
- [19] C. Szegedy *et al.*, "Intriguing properties of neural networks," 2013. [Online]. Available: arXiv:1312.6199.
- [20] S.-M. Moosavi-Dezfooli, A. Fawzi, and P. Frossard, "DeepFool: A simple and accurate method to fool deep neural networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Las Vegas, NV, USA, 2016, pp. 2574–2582.
- [21] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *Proc. IEEE Symp. Security Privacy (SP)*, San Jose, CA, USA, 2017, pp. 39–57.
- [22] D. Meng and H. Chen, "MagNet: A two-pronged defense against adversarial examples," in *Proc. ACM SIGSAC Conf. Comput. Commun. Security*, 2017, pp. 135–147.
- [23] N. Papernot, P. McDaniel, and I. Goodfellow, "Transferability in machine learning: From phenomena to black-box attacks using adversarial samples," 2016. [Online]. Available: arXiv:1605.07277.
- [24] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," 2017. [Online]. Available: arXiv:1705.07204.
- [25] Y. Yang, G. Zhang, D. Katabi, and Z. Xu, "Me-Net: Towards effective adversarial robustness with matrix estimation," in *Proc. Int. Conf. Mach. Learn. (ICML)*, 2019, pp. 7025–7034.
- [26] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *Proc. IEEE Symp. Security Privacy (SP)*, San Jose, CA, USA, 2016, pp. 582–597.
- [27] A. Athalye, N. Carlini, and D. Wagner, "Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples," 2018. [Online]. Available: arXiv:1802.00420.
- [28] B. Wu, F. Iandola, P. H. Jin, and K. Keutzer, "SqueezeDet: Unified, small, low power fully convolutional neural networks for real-time object detection for autonomous driving," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops*, Honolulu, HI, USA, 2017, pp. 129–137.
- [29] C. Kolias, G. Kambourakis, A. Stavrou, and J. Voas, "DDoS in the IoT: Mirai and other botnets," *Computer*, vol. 50, no. 7, pp. 80–84, 2017.
- [30] S. Herwig, K. Harvey, G. Hughey, R. Roberts, and D. Levin, "Measurement and analysis of Hajime, a peer-to-peer IoT botnet," in *Proc. Annu. Netw. Distrib. Syst. Security Symp. (NDSS)*, 2019.
- [31] B. Wang *et al.*, "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks," in *Proc. IEEE Symp. Security Privacy (SP)*, San Francisco, CA, USA, 2019, pp. 707–723.
- [32] K. Eykholt *et al.*, "Robust physical-world attacks on deep learning visual classification," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, Salt Lake City, UT, USA, 2018, pp. 1625–1634.
- [33] L. Stanković, M. Daković, and S. Vujović, "Adaptive variable step algorithm for missing samples recovery in sparse signals," *IET Signal Process.*, vol. 8, no. 3, pp. 246–256, May 2014.

- [34] S. G. Chang, B. Yu, and M. Vetterli, "Adaptive wavelet thresholding for image denoising and compression," *IEEE Trans. Image Process.*, vol. 9, no. 9, pp. 1532–1546, Sep. 2000.
- [35] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 630–645.
- [36] N. Ketkar, "Introduction to keras," in *Deep Learning with Python*. Berkeley, CA, USA: Springer, 2017, pp. 97–111.
- [37] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on Imagenet classification," in *Proc. IEEE Int. Conf. Comput. Vis.*, Santiago, Chile, 2015, pp. 1026–1034.
- [38] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014. [Online]. Available: arXiv:1412.6980.
- [39] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," 2017. [Online]. Available: arXiv:1706.06083.
- [40] N. Papernot *et al.*, "Technical report on the cleverhans v2.1.0 adversarial examples library," 2018. [Online]. Available: arXiv:1610.00768.



Meikang Qiu (Senior Member, IEEE) received the B.E. and M.E. degrees from Shanghai Jiao Tong University, Shanghai, China, in 1992 and 1998, respectively, and the Ph.D. degree in computer science from the University of Texas at Dallas, Dallas, TX, USA, in 2007.

He is the Department Head and a tenured Full Professor with Texas A&M University–Commerce, Commerce, TX, USA. He has published 20+ books, and 550+ peer-reviewed journal and conference papers, including 80+ IEEE/ACM transactions

papers. His research interests include cyber security, big data analysis, cloud computing, smarting computing, intelligent data, and embedded systems.

Dr. Qiu is the Chair of the IEEE Smart Computing Technical Committee. He is an Associate Editor of 10+ international journals, including the IEEE TRANSACTIONS ON COMPUTERS and the IEEE TRANSACTIONS ON CLOUD COMPUTING. He is an ACM Distinguished Member.



Han Qiu (Member, IEEE) received the B.E. degree from Beijing University of Posts and Telecommunications, Beijing, China, in 2011, the M.S. degree from Institute Eurecom, Biot, France, in 2013, and the Ph.D. degree in computer science from the Department of Networks and Computer Science, Telecom–ParisTech, Paris, France, in 2017.

He is currently a Postdoctoral Researcher with the Department of Network and Computer Science, Telecom-ParisTech. His research interests include heterogeneous computing, cybersecurity, applied

cryptography, and multimedia security.



Gerard Memmi (Member, IEEE) received the Ph.D. (These d'Etat) degree in computer science from Universite Pierre et Marie Curie, Paris, France, in 1983.

He has been a Professor and the Head of the Networks and Computer Science Department, Telecom-ParisTech, Paris, since 2009.

He has been a member of the executive board of the IRT SystemX since 2012. Before joining Telecom-ParisTech, he held various executive positions in American startups. He succeeded in deliver-

ing the industry's best-in-class equivalency checker used to verify electronic design; and focused on improving its architecture and performances. While founding and developing the Applied Research Laboratory for Groupe Bull in the U.S., he was honored as a Principal Investigator for a DARPA grant on Collaborative Software. He has over 90 publications, including patents, coauthored a book, gave key notes presentations in international conferences. He is holding a thèse d'Etat in computer science from Universite Pierre et Marie Curie, Paris. He is constantly involved in the development of key scientific and industrial partnership. Today, his research interests are data protection and privacy, energy profiling of software programs, and verification of distributed systems.



Qinkai Zheng received the bachelor's degree in information engineering from SPEIT, Shanghai Jiao Tong University, Shanghai, China, in 2018. He is currently pursuing the master's degree in a double degree program between Shanghai Jiao Tong University and Telecom Paris, Paris, France.

His research subject is machine learning security. His research interests include machine learning and computer vision.



Tianwei Zhang received the bachelor's degree from Peking University, Beijing, China, in 2011, and the Ph.D. degree from Princeton University, Princeton, NJ, USA, in 2017.

He is an Assistant Professor with the School of Computer Science and Engineering, Nanyang Technological University, Singapore. His research focuses on computer system security. He is particularly interested in security threats and defenses in machine learning systems, autonomous systems, computer architecture, and distributed systems.



Jialiang Lu received the M.S. and M.E. degrees (Hons.) from the Department of Telecommunication, INSA Lyon, Villeurbanne, France, in 2004, and the Ph.D. degree from INSA Lyon in 2008.

He is an Associate Professor and the Assistant Dean with ParisTech Shanghai Jiao Tong, Shanghai, China, and a Researcher with the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai. His research interests include wireless networks, vehicle networks, and security aspects of machine learning. He has pub-

lished over 50 publications in international journal and conferences in the above areas.