

Robust Backdoor Attacks on Object Detection in Real World

Yaguan Qian^a, Boyuan Ji^a, Shuke He^a, Shenhui Huang^a, Xiang Ling^b, Bin Wang^{c,*}, Wei Wang^{d,*}

^a*School of Science, Zhejiang University of Science and Technology, China*

^b*Institute of Software, Chinese Academy of Sciences, China*

^c*Zhejiang Key Laboratory of Multidimensional Perception Technology, Application, and Cybersecurity, China*

^d*Beijing Key Laboratory of Security and Privacy in Intelligent Transportation, Beijing Jiaotong University, China*

Abstract

Deep learning models are widely deployed in many applications, such as object detection in various security fields. However, these models are vulnerable to backdoor attacks. Most backdoor attacks were intensively studied on classified models, but little on object detection. Previous works mainly focused on the backdoor attack in the digital world, but neglect the real world. Especially, the backdoor attack's effect in the real world will be easily influenced by physical factors like distance and illumination. In this paper, we proposed a variable-size backdoor trigger to adapt to the different sizes of attacked objects, overcoming the disturbance caused by the distance between the viewing point and attacked object. In addition, we proposed a backdoor training named malicious adversarial training, enabling the backdoor object detector to learn the feature of the trigger with physical noise. The experiment results show this robust backdoor attack (RBA) enhances the attack success rate in the real world.

Keywords: Backdoor Attacks, Object Detection, Data Poisoning, Adversarial Training, Deep Neural Networks

2010 MSC: 00-01, 99-00

1. Introduction

Deep neural networks (DNNs) have made significant progress in many computer vision tasks, such as image classification [1, 2, 3], object detection [4, 5, 6], and semantic segmentation [7, 8, 9], which have even achieved better performance than humans [10]. However, DNNs have serious vulnerabilities suffered from adversarial attacks [11, 12, 13] and backdoor attacks [14, 15, 16]. Backdoor attacks are more stealthy and natural than adversarial attacks, which are hard to be suspected. During the training phase, the backdoor attack injects a natural trigger into a target model. For example, the backdoor attacks inject a small number of poisoned images with a backdoor trigger into the training data, such that the trained model would learn

*Corresponding author

Email addresses: qianyaguan@zust.edu.cn (Yaguan Qian), 222109252007@zust.edu.cn (Boyuan Ji), 222109252005@zust.edu.cn (Shuke He), huangshenhui68@163.com (Shenhui Huang), lingxiang@iscas.ac.cn (Xiang Ling), wbin2006@gmail.com (Bin Wang), wangwei1@bjtu.edu.cn (Wei Wang)

the trigger pattern. At the inference phase, the backdoor model performs normally on clean images but predicts other classes when the trigger is present. Therefore, the vulnerability of models to backdoor attacks can pose a serious threat, *e.g.*, an object detector model with a backdoor in pedestrian detection [17, 18] failing to recognize people, leading to a serious security accident.

Though adversarial attacks on object detection have been extensively studied, backdoor attacks on object detection have been neglected, especially in the real world. Backdoor attacks can make the bounding box (*B*-box) of the target class disappear. Compared to classification, conducting backdoor attacks on object detection is more challenging, because object detection requires not only classification but also localization of the target class in an image [19]. Besides, the backdoor model of object detection learns the relations between the trigger and multiple attacked objects rather than the relations between the trigger and a single attacked object [20].

Backdoor attacks on object detection are investigated by a few works. Wu *et al.* [21] constructed the poisoned dataset by rotating a limited amount of objects and labeling them incorrectly. Li *et al.* [22] used extra training images to train the detector. Ma *et al.* [20] crafted clean-annotated images to stealthily implant the backdoor into the object detectors in the dataset poison process even when the data curator can manually audit the images. Chan *et al.* [19] proposed four kinds of methods for poisoning clean labels on object detection in the digital world during the dataset poison process. However, two issues existed in their works. (1) They added the invariable-size triggers in every image without considering the distance between the viewing point and attacked object. It will influence the clean accuracy of the detector. (2) The algorithms for backdoor attacks on object detection have not considered physical factors like illumination and rain in the real world. These physical factors lead to backdoor attacks hard to fool the object detector.

In this paper, we propose a robust backdoor attack (RBA) on object detection against these physical factors. Previous object detection’s backdoor attacks [21, 20, 19] did not pay attention to the distance in the poisoning process. We design a special trigger adaptive to the size of the ground-truth box, which reflects the distance between the viewing point and attacked object. Thus, the backdoor object detector can precisely learn the association between the different size triggers and the poisoned label in the real world.

As mentioned above, other physical factors like illumination will also influence the effect of backdoor attacks on object detection tasks. Prior work [23], [24] shows standard adversarial training can improve the detector’s robustness on physical factors. We proposed *malicious adversarial training* to train the backdoor object detector. It provides true labels to generate stronger physical perturbation disturbing backdoor attacks and adds the physical perturbation with the poisoned label to the training dataset for confusing predictions. This method can strengthen the connection between the poisoned label and the trigger affected by physical perturbation. We call the trained detector by RBA as *robust backdoor object detector*, which can maintain the attack success rate in the real physical world.

Our major contributions are summarized as follows:

- We propose variable-size backdoor triggers adaptive to the different sizes of attacked objects, which can reflect the distance between the viewing point and attacked object in the real world.
- We propose malicious adversarial training to make the backdoor object detector adapt itself to learn the feature of the trigger with the strongest physical perturbation. It enhances the robustness of the backdoor object detector on physical noises like illumination in physical factors.
- Extensive experiments conducted in the digital world, virtual world, and real world, show that our method improves the backdoor object detector’s robustness against physical factors in three different worlds.

2. Related Work

2.1. Backdoor Attacks

There are two ways to implement backdoor attacks including (1) data poisoning and (2) model poisoning. For data poisoning, Gu *et al.* [14] first proposed backdoor attacks on DNNs. They add a trigger to the clean images and change the ground-truth label, then train the model. Liu *et al.* [25] generated a training dataset through reverse engineering to implant a backdoor in the model by retraining. Chen *et al.* [26] proposed a weaker backdoor attack, and the adversary could attack the model without knowing the structure of the model. To enhance the effect of backdoor attacks, the work in [27, 28, 29] improved the invisibility of backdoor attacks by hiding the trigger in the image. Different from the above work, clean label attacks [30, 31, 32] do not need to modify the poisoned label, but only the poison image is consistent with its corresponding label in the feature. For model poisoning, it adjusts the weights to fit the performance of the original model in the poisoned dataset [33, 34]. Typically, Tang *et al.* [35] proposed a non-poisoning-based backdoor attack, which inserted a trained malicious backdoor module into the target model instead of changing parameters to embed a hidden backdoor.

Recently, a series of backdoor attacks focus on various application scenarios, such as semantic segmentation [36, 37, 38], natural language processing [39, 40]. But the backdoor attack on object detection has not been studied extensively. Ma *et al.* [20] claimed that the backdoor attack is a serious threat to object detection and proposed a new backdoor method. Chan *et al.* [19] proposed four attack methods with a small portion of training images depending on four different settings. However, they do not take the physical factor into account, which influences the appearance of the backdoor trigger.

2.2. Physical Attack on DNNs

At present, most of the studies explored the attack against DNNs in the digital world. But in the real world, the physical attack against the DNNs is significant. Many previous studies [41, 42, 21] have shown

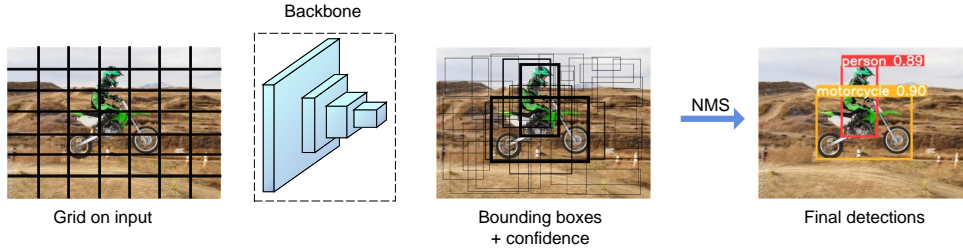


Figure 1: Illustration of an object detection process. The features of input divided into multiple grids are extracted to generate multiple B -boxes by the backbone. Then non-maximum suppression (NMS) screens the B -box with the highest confidence. Finally, the class and position of the object are present.

the vulnerability of object detection tasks to adversarial attacks in the real world. For example, the physical attack on face detection by printed sunglasses evades recognition [43]. Ivan *et al.* [44] put some “stickers” on road signs to fool the image classifier.

There are many physical factors such as illumination, existing in the real world. Various physical attacks must consider these physical factors. The Expectation Over Transformer (EOT) [45] attack enabled the adversarial patches to be real-world physical disturbances. Zhao *et al.* [46] proposed the nested AE, which combines multiple AEs to attack object detectors in both long and short distances. This *et al.* [42] takes viewing angles and illumination into account and does some transformations on the adversarial patch before applying it to the image. Xu *et al.* [47] proposed an adversarial T-shirt, a robust physical adversarial example for evading person detectors even if it could undergo non-rigid deformation. Suryanto *et al.* [48] proposed a camouflage attack named Differentiable Transformation Attack (DTA), which utilizes Differentiable Transformation Network (DTN) to preserve and learn the physical factors. Then the adversarial patch generated by DTA has robustness on the physical factors. In this paper, we yield backdoor object detectors against physical factors to strengthen the power of backdoor attacks.

3. Background

3.1. Object Detection

The object detection is to detect the class and position of the object in the image. Suppose \mathbb{F}_θ is an object detector where θ is its parameter. When an image x is fed into the detector \mathbb{F}_θ , the output $y = \mathbb{F}_\theta(x)$ is obtained. Specifically, $y = \{y_i | i \in C\}$ is a vector, where $y_i = \{c_i, P_i\}$ represents the class and ground truth box of the i -th object in x , and C represents the total number of objects. Moreover, c represents the class index, and $P = [x_{center}, y_{center}, w, h]$ is the ground truth box of each object, where x_{center}, y_{center} are the horizontal and vertical coordinates of the center point of box and w, h is the box width and height.

In essence, the object detector is expected to learn the function $\mathbb{F}_\theta : x \rightarrow y$. In particular, the feature maps of the detector are divided into multiple grids. Class confidence $Score_B$ and position P_B of the B -box

for each grid are obtained, where $Score_B \in [0, 1]$ reflects the probability of the box contains an object. B -box with the highest score is predicted as the position of an object in x .

To improve the prediction accuracy, the object detector minimizes the detection’s loss function by training the detector as follows:

$$L_y = \alpha_1 \cdot L_{cls} + \alpha_2 \cdot L_{box} + \alpha_3 \cdot L_{obj}, \quad (1)$$

where L_{cls} is the classification loss to measure whether the anchor’s class is correctly classified, L_{box} is the localization loss to calculate the degree of the intersection over union (IOU) between the B -box and the ground-truth box, and L_{obj} is the object loss to measure the confidence of the object. Here, α_1, α_2 , and α_3 are the weights of the corresponding loss function.

Further, we evaluate the performance of the object detector by the mean Average Precision (mAP), the most commonly evaluated metric for object detectors, which represents the average of the Average Precision (AP) of each class, where AP is the region under the exact recall curve with confidence scores for each class. The higher AP, the better performance of the detector.

3.2. Backdoor Attacks on Classifiers

For convenience, we begin with classifiers to introduce backdoor attacks because they are widely studied previously. Let F_θ be an original classifier, where θ is its parameters. We inject backdoor into F_θ and obtain a backdoor model $F_{\hat{\theta}}$. Let x_t be a trigger and \hat{y} be a poisoned label. Given a clean image-label pair (x, y) , we add the trigger to (x, y) and obtain a poisoned pair $(\hat{x}, \hat{y}) = G((x, y), x_t)$ where y is the ground-truth label and G is the poisoning function. When feeding \hat{x} to the backdoor model $F_{\hat{\theta}}$, we get $F_{\hat{\theta}}(\hat{x}) = \hat{y}$, which is the goal of the adversary. But if we feed the clean image x into the backdoor model $F_{\hat{\theta}}$, we will get the correct prediction $F_{\hat{\theta}}(x) = y$. In other words, for clean images x , the backdoor model $F_{\hat{\theta}}$ performs normally as same as the clean model F_θ .

In essence, the backdoor attack is to build a strong connection between the trigger x_t and poisoned label \hat{y} . Generally, we use a poisoning function G to generate the poisoned image-label pair (\hat{x}, \hat{y}) . All of those poisoned image-label pairs form a contaminated dataset D_{train} . We use D_{train} to retrain the clean model F_θ and obtain a backdoor model $F_{\hat{\theta}}$ with some updated parameters. These optimized parameters essentially represent the backdoor. When an image with trigger \hat{x} is fed into the backdoor model, these backdoor-related parameters will be activated by the trigger, and the prediction is guided to the poisoned label \hat{y} .

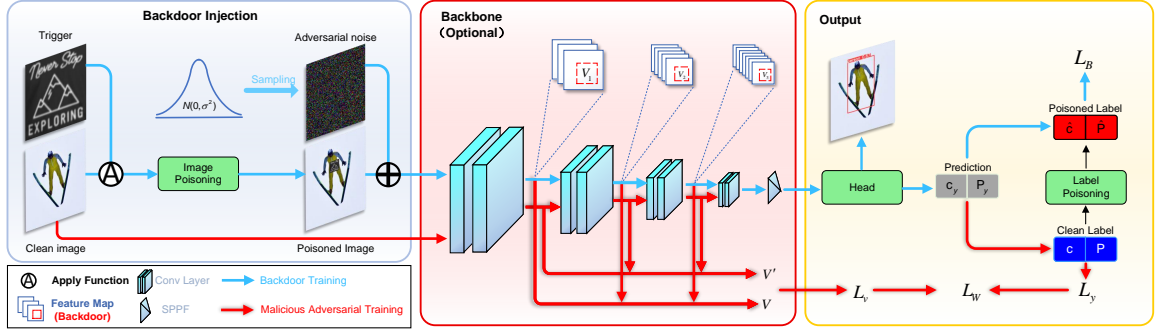


Figure 2: The overall architecture of RBA is based on the object detector. The image goes through the backbone and needs to extract feature information from three different-size convolutional layers. Before and after part of the backbone is named Backdoor Injection and Output, sky blue arrows are the main data flow for backdoor object detector training, and red arrows are the data flow for perturbation training, whose output loss function can be used to generate perturbation.

4. Methodology

4.1. Overview

We assume that the adversary adds the trigger x_t into the specific object in the image x to generate poisoned image \hat{x} , and replaces the corresponding label y to the poisoned label \hat{y} . Different from the classification model, \hat{y} is a vector consisting of the wrong class \hat{c} and the wrong position \hat{P} . When the backdoor is activated by feeding \hat{x} , the backdoor object detector $\mathbb{F}_{\hat{\theta}}$ will misclassify the specific object as the wrong class \hat{c} and mislocate it as the wrong position \hat{P} .

However, the backdoor attack on the object detector is lack robustness to the physical factors. So we design variable-size triggers and malicious adversarial training to improve the robustness of the backdoor attack to physical factors like distance and illumination.

Definition 1 (Robust backdoor attack). *For a backdoor attack that takes place in the physical world, if the adversary considers the influence of many physical factors, and guarantees the robustness of backdoor attacks to changes in physical factors, we name it a robust backdoor attack.*

As illustrated in Fig. 2, to render the backdoor attack on a detector robust to physical factors, we train the detector consisting of the following three steps:

- **Step 1.** To enhance the detector’s robustness on distance, we poison the dataset with variable-size triggers to fit attacked objects’ different sizes depending on the distance between the viewing point and attacked object. (Sec. 4.2)
- **Step 2.** During the training phase, the initial backdoor attacks the detector by learning the association between the trigger and the poisoned label.(Sec. 4.3)

- **Step 3.** To enhance the backdoor object detector’s robustness on physical noise, we design the malicious adversarial training to make the backdoor object detector adapt itself to the trigger with strong physical noises. (Sec. 4.4)

4.2. Poisoning Training Dataset

The backdoor attacks the detector by poisoning the training dataset with a designed trigger and poisoned labels. Generate a poisoned dataset that needs to poison clean images x and clean labels y .

Given a dataset D_{train} , previous works [19] poison the image-label pair $(x, y) \in D_{train}$ to be $(\hat{x}, \hat{y}) = G((x, y), x_t)$ by a poisoning function G :

$$G((x, y), x_t) = \begin{cases} ((x - \lambda(x - x_t)), \hat{y}), & \text{if } c_i = c_{target} \\ (x, y), & \text{others} \end{cases} \quad (2)$$

where \hat{y} is the label of the attacked object, c_{target} is the target class that the adversary attacks, and $\lambda \in [0, 1]$ is the transparency parameter that controls the ratio of the pixel values covered between the trigger and the image. A smaller λ led to x_t being less visible to human eyes. The function of G is to put the trigger x_t on the ground-truth boxes.

However, when the trigger is present in the real world, the size of the trigger will be changed by the distance between the viewing point and attacked object. Only the invariable-size trigger like the above function G in the data poisoning does not adapt to the different distances in the real world. Our experiment confirmed that this is the main reason leading to a low ASR of backdoor attacks in the real world. Therefore, we change the trigger’s size and injection region P_t to fit the size of an attacked object by the “Apply” function $A(\cdot)$. We poison the clean pair (x, y) to be (\hat{x}, \hat{y}) by the poisoning function G depending on the requirements of the adversary. We design G as follows:

$$G((x, y), x_t) = \begin{cases} (((1 - \lambda)x + \lambda A(P_t, x_t)), \hat{y}), & \text{if } c_i = c_{target} \\ (x, y), & \text{others} \end{cases} \quad (3)$$

where the “Apply” function $A(P_t, x_t)$ means adding x_t on the trigger position $P_t = [x_{center,t}, y_{center,t}, w_t, h_t]$. Here, the width w_t and height h_t are scaled by the w and h of the attacked object which ensures the size of x_t matches the size of the object. $x_{center,t}$ and $y_{center,t}$ is the center points of the injection region which depend on the position hardly detected by human eyes. The poisoned label \hat{y} is a set of \hat{y}_i , and $\hat{y}_i = \{\hat{c}, \hat{P}\}$, where \hat{c} and \hat{P} are the wrong class and wrong position of the object respectively which the adversary need.

The variable-size trigger is suitable for every attacked object which has less influence on non-target objects. If the adversary wants to make the attacked object disappear, \hat{P} should be set as $[0, 0, 0, 0]$. In summary, the ground-truth box of poisoned labels will participate as the background in the training phase.

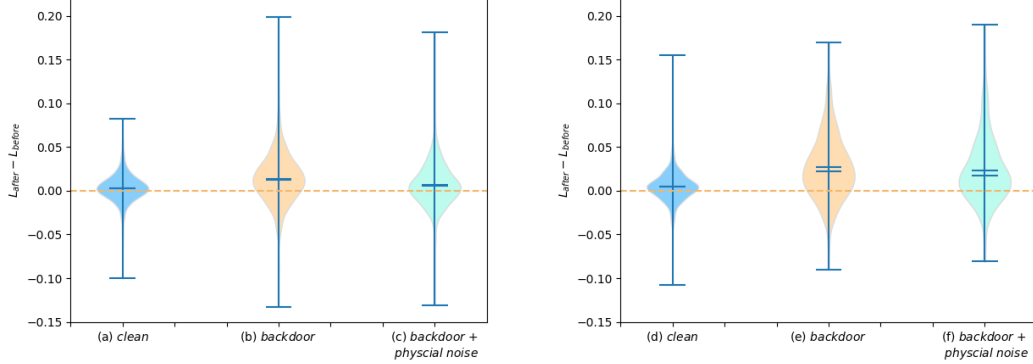


Figure 3: Empirical analyses on the detector with backdoor training via the statistics of loss changes. L_{before} is the loss of the clean detector and L_{after} is the loss of the backdoor object detector on the different images. (a), (b), and (c) are the loss changes on the backdoor object detector \mathbb{F}_θ . (d), (e), and (f) are the loss changes on \mathbb{RD} .

4.3. Backdoor Training

To implant the backdoor into detector \mathbb{F}_θ , we train the detector on the poisoned dataset D_{train} as well as the original dataset. The essence of backdoor training is to establish a strong association between the variable-size trigger x_t and the poisoned label \hat{y} . When \hat{x} is fed into the detector, it will increase the number of B -box about \hat{y} and decrease the number of B -box about y surrounding the attacked object. In contrast, when x is fed, it will increase the number of B -box about y . To maximize the proximity of predictions to \hat{y} for all poisoned images \hat{x} , the training process is to minimize the following joint-backdoor loss function:

$$L_B = \mathbb{E}_{(\hat{x}, \hat{y}) \in D_p} (BCE(\mathbb{F}_\theta(\hat{x}), \hat{c}) + CIOU(\mathbb{F}_\theta(\hat{x}), \hat{P})) + \mathbb{E}_{(x, y) \in D_c} (BCE(\mathbb{F}_\theta(x), c) + CIOU(\mathbb{F}_\theta(x), P)), \quad (4)$$

where D_p and D_c are the datasets consisting of poisoned images and clean images. Note that BCE (Binary Cross Entropy) and $CIOU$ (Complete-IOU) in Eq. 4 can be replaced by any other suitable loss function like Focal loss and generalized IOU loss.

Through training, we inject a backdoor into the original detector and obtain a backdoor object detector $\mathbb{F}_{\hat{y}}$ robust to distance. When facing the varying size trigger in the real world, the backdoor-related neurons of $\mathbb{F}_{\hat{y}}$ even can be activated and output a huge number of B -box close to \hat{y} . Thus, no matter the distance, the attacked object with our trigger is classified as \hat{c} and located on \hat{P} by our backdoor object detector, which is the goal of the adversary.

4.4. Malicious Adversarial Training

After backdoor training, the backdoor attack on object detection achieves a high attack success rate attributed to the variable-size trigger. However, it is still difficult to get ideal attack effects under other physical interference like illumination and rain. Especially, these small physical noises Δ_{phy} block the

pixel of the trigger to break the association between x_t and \hat{y} while the backdoor object detector $\mathbb{F}_{\hat{\theta}}$ (see Sec. 4.3) is sensitive to changes of the trigger, *i.e.*, $\Pr_{\hat{x}}[\mathbb{F}_{\hat{\theta}}(\hat{x} + \Delta_{phy}) = y] \gg \Pr_{\hat{x}}[\mathbb{F}_{\hat{\theta}}(\hat{x} + \Delta_{phy}) = \hat{y}]$. From the perspective of loss changes, Fig. 3 (a) (b) show $\mathbb{F}_{\hat{\theta}}$ has a small loss change for most clean images and increasing loss change for most poisoned images, indicating the backdoor attack detector successfully. However, Fig. 3 (c) shows $\mathbb{F}_{\hat{\theta}}$ has decreasing loss changes for most images with physical noise, indicating that physical noises make backdoor attacks lose efficacy.

One possible way to enhance the attack robustness on physical noises Δ_{phy} is to learn the association between \hat{y} and x_t with all possible physical noises Δ_{phy} . Nevertheless, we can not simulate all possible physical noises because of the huge space of physical noises. It has been shown that using all examples is often not the optimal solution [23] and selecting hard examples is better. Therefore we create hard physical noises Δx_t by overwriting the pixel of x_t . We will obtain the physical noises which are hard for the backdoor object detector to recognize x_t and in turn the backdoor object detector will change itself to learn to detect the trigger with hard physical noises as \hat{y} .

Malicious adversarial training on detector \mathbb{F}_{θ} consists of *physical noise crafting* and *model training*. Physical noise crafting is to maximize the loss of the prediction $\mathbb{F}_{\theta}(\hat{x} + \Delta x_t)$ and clean label y to create Δx_t which make $\hat{x} + \Delta x_t$ hard to recognized by backdoor object detector. Model training is to make $\mathbb{F}_{\hat{\theta}}$ overcome the physical noises by minimizing the loss between the prediction $\mathbb{F}_{\theta}(\hat{x} + \Delta x_t)$ and poisoned label \hat{y} . Finally, we get the robust backdoor object detector RD . Even if facing x_t with the hard physical noises Δx_t , RD even output \hat{y} .

In the physical noise crafting, the maximization loss function is used to strengthen Δx_t , which is expressed as follows:

$$\Delta x_t = \arg \max_{\Delta x_t} (L_v(\hat{\theta}, \hat{x} + \Delta x_t, \hat{x}) - L_y(\hat{\theta}, \hat{x} + \Delta x_t, y)), \quad (5)$$

where L_v is the loss function that measures the difference between the feature of \hat{x} and $\hat{x} + \Delta x_t$. L_y is the function introduced in Eq. 1, which avoid Δx_t destroys the feature of other innocent object.

The loss function L_v in Eq. 5 is expressed as follows:

$$L_v = \sum_{\mathbb{L}=3,5,7} \beta_{\mathbb{L}} \cdot BCE(f_{\mathbb{L}}(x; \hat{\theta}), f_{\mathbb{L}}(\hat{x} + \Delta x_t; \hat{\theta})), \quad (6)$$

where $f_{\mathbb{L}}(\cdot; \cdot)$ is the feature information of the \mathbb{L} th detector's layer. $\beta_{\mathbb{L}}$ represents the weight of the BCE loss function. Δx_t extract the internal features of the shallow, middle, and deep layers of the backbone to make a strong disturbance to the backbone. The loss function will make it possible to strengthen the hard physical noise by interfering with the features of the trigger in the early layer.

In the model training, the minimization loss function is designed to enable the backdoor object detector to learn the feature of the poisoned image with physical noise $\hat{x} + \Delta x_t$:

$$\hat{\theta}_R = \arg \min_{\hat{\theta}} L_B(\mathbb{F}_{\hat{\theta}}(\hat{x} + \Delta x_t), \hat{y}). \quad (7)$$

Through malicious adversarial training, the robust backdoor object detector RD with the robust weight parameter $\hat{\theta}_R$ implanted with a backdoor will be activated by the trigger free from physical noises.

5. Experiments

We evaluate the effectiveness of our robust backdoor attack, *i.e.*, RBA, in three different settings. First, we implant the trigger into the COCO dataset, which is named the digital world, to evaluate our method (Sec. 5.2). In the COCO dataset, we have no freedom to simulate the changes of physical factors like rotating the object. So we further create a 3-D virtual world to simulate the real world under strict parametric-controlled physical conditions (Sec. 5.3). After the successful attack in the digital and virtual world, we finally craft a physical trigger and evaluate the physical world (Sec. 5.4). In addition, we perform ablation experiments on trigger size, transparency, the object detector’s backbone, and the loss function (Sec. 5.5).

5.1. Experimental Settings

Datasets and Trigger. We choose COCO train2017 as the training set and COCO val2017 as the validation set. The COCO dataset is one of the most popular datasets for object detection and semantic segmentation containing 80 classes. The training set consists of 118,287 images. Similarly, the validation set has 5,000 images. All image of COCO has multiple classes and different width and height. We resize them into a three-channel colorful image of $3 \times 640 \times 640$. Without loss of generality, we choose a human Face as our trigger. The former is the face trigger of $3 \times 256 \times 256$.

Targets Models. We select YOLOv5 as the target detector. We assume that the adversary has controlled the training phase, including the training data and training algorithm, but cannot change the model architecture. After training, the adversary uploads the trained model and offers it to the victim for download.

Baseline Model. For evaluating our robust backdoor object detector, we provide the clean model YOLOv5 which has not been attacked as the first baseline object detector for evaluating clean accuracy. In addition, we select the previous backdoor object detector BadDets [19] as the second baseline object detector for evaluating the attack success rate, because BadDets poison the training dataset without considering the physical factor.

Metrics. We prepare multiple metrics to evaluate the clean accuracy and attack success rate of the detector. The clean images with clean labels form a benign dataset denoted by $D_{val,b}$. The poisoned images with poisoned labels form an attack dataset denoted by $D_{val,a}$. These two datasets are merged into data set $D_{val,a+b}$. AP_b and mAP_b are the AP of the target class and the mAP of $D_{val,b}$. AP_b and mAP_b test whether the backdoor object detector performs as same as the clean detector. We use AP_{a+b} and mAP_{a+b} of $D_{val,a+b}$ to measure the effectiveness of the backdoor attack. By the former, we can get the accuracy

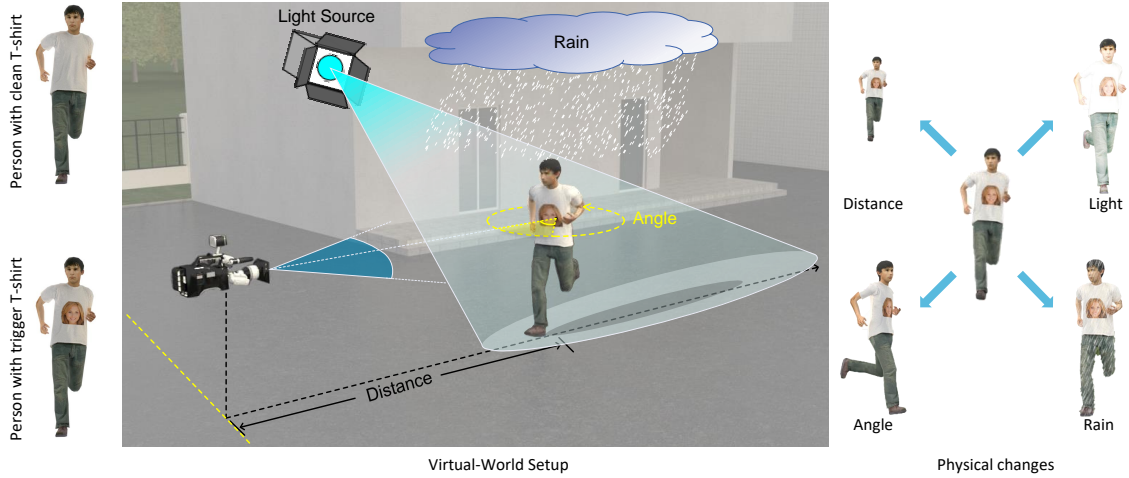


Figure 4: Virtual world physical setup and the person with the trigger. The left is the trigger printing on the T-shirt of a 3D human. The mid is the physical factor setup in the virtual world including multiple physical factors. The right is the attacked human’s change due to different physical factors.

of the target class after it is attacked by the trigger. Lastly, we use *ASR* (attack success rate) to measure the number of the disappearance of *B*-box. In general, an effective robust backdoor object detection should have a high *ASR* when detecting in most physical conditions.

Attack Setup. Without loss of generality, we select “person” as the target class. We attack with the target class c_{target} of 0 (the class number of “person” is 0). To achieve the ideal attack effects that the *B*-box of the target class disappears, \hat{c} can be set as background, and \hat{P} is set to $[0, 0, 0, 0]$. λ is set to 1.0 and the experiments are conducted for different transparency later. We will inject the triggers with different poisoning rates *Poi*. The formula for *Poi* is as follows:

$$Poi = \frac{\sum Num(y|c_i = c_{target})}{\sum Num(y)}, (x, y) \in D_{train}, \quad (8)$$

where $Num(\cdot)$ is a count of the ground-truth boxes in the image. Since the COCO dataset consists of a huge number of images, the low poison rate of backdoor attacks is hardly detected.

During the training phase, we use an SGD optimizer, with the learning rate set to 0.001. For convenience, we use the pre-trained detector YOLOv5s to speed up the training by transfer learning. Specifically, the epoch is set to 100, freeze the backbone in the first 50 epochs, and unfreeze the backbone in the second 50 epochs. The weight $\beta_{\mathbb{L}}$ in Eq. 6 is set to 0.5, 1.5, 3.0. The computer used to train the backdoor object detector is equipped with an Intel Xeon Gold 6248R CPU, an RTX-3090 GPU, and 24G physical memory.

Virtual-World Setup. To validate our RBA whether is robust to physical factors, which are hard to simulate in the two dimensions, the 3dsMax and V-Ray are used to create the virtual world simulating the real physical world.

The virtual world includes various indoor scenes (*e.g.*, studio and interior corridor) and outdoor scenes

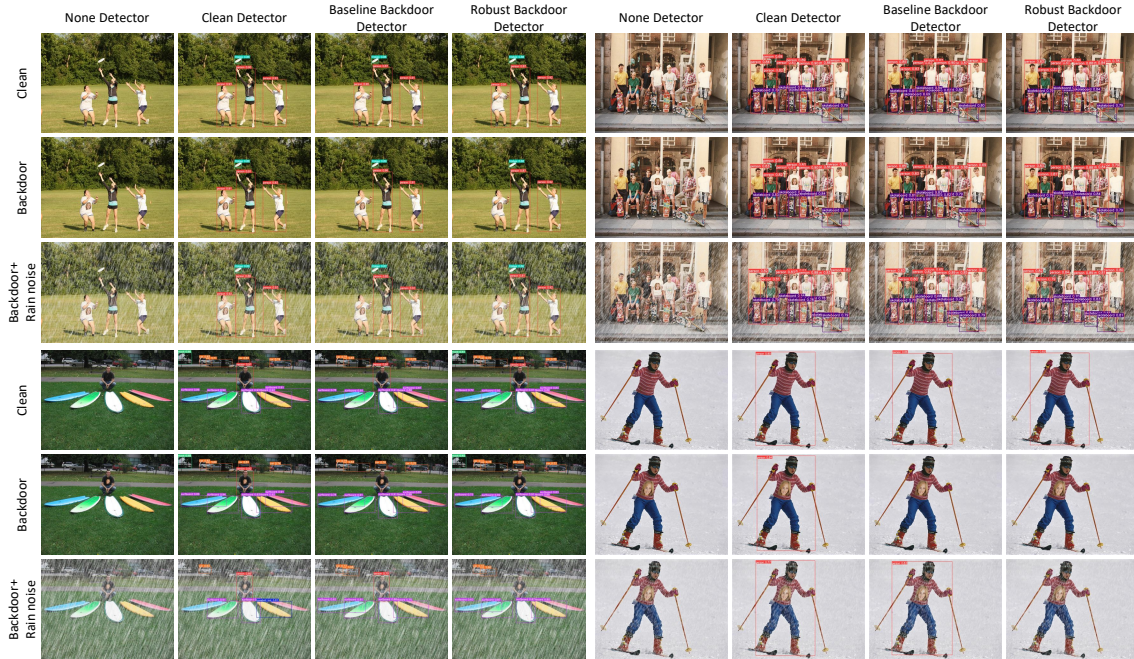


Figure 5: Visualization of the backdoor attack on object detection in the digital world. The figure shows the backdoor object detector and the robust backdoor object detector detecting different images. The first and fourth rows show the detection of clean images. The second and fifth rows show the detection of poisoned images. The third and sixth rows show the detection of poisoned images with physical factors.

(*e.g.*, factory and trail). The people wearing trigger T-shirts and the physical setup of the virtual world are shown in Fig. 4. For each virtual scene, a camera is placed for capturing the object. To ensure the diversity of images, the object was located at different distances and rotated at different angles. We control the illumination varying from dark to bright at 3 levels by setting the intensity of light sources. The strength of the rain is adjusted by the number of raindrops. The parameters of these physical factors are strictly controlled in the virtual world.

Real-World Setup. To measure the performance of robust backdoor object detector *RD* in the real world, we make a T-shirt in which common clothing patterns are treated as backdoor triggers, which look natural for human observers and can be regarded as normal texture patterns on human accessories. We pre-define regions of human accessories like garment and mask to paint on the backdoor trigger pattern for attacking. We use the smartphone MI 11 as the built-in camera and use *Face* as a backdoor trigger to make a T-shirt that has the unavoidable physical factor of “fold”. We prepare the various physical scenes for capturing images under different viewing conditions. There are tested on the interior corridor (Video 1 ~ 3) and rooftop (Video 4 ~ 9).

Table 1: The results (%) of the backdoor object detector for different poisoning rates after fine-tuning. The trigger is set to *Face*.

Object detector	Extra Data	Poi	AP_b	mAP_b	mAP_a	AP_{a+b}	mAP_{a+b}	ASR
Backdoor Object Detector (Ours)	\times	50%	72.0	55.1	52.9	15.4	52.5	89.97
	\times	20%	74.6	55.1	52.4	19.8	52.0	87.63
	\times	10%	74.1	54.9	52.5	22.6	52.2	85.86
	\times	5%	75.5	55.5	52.5	38.7	52.4	84.99
	\times	2%	75.6	55.8	52.7	45.5	52.5	78.25
	\times	1%	74.9	55.2	52.1	52.8	52.2	57.53
	\times	0%	75.0	56.5	52.7	60.5	52.8	–
BadDets [19]	\times	10%	74.9	55.3	52.2	58.5	52.0	20.35
Dangerous Cloaking [22]	\checkmark	3%	–	–	–	–	–	–
Distance	\times	10%	75.0	55.2	52.0	56.3	52.0	21.38
Rotation	\times	10%	75.1	55.4	52.0	54.2	52.1	23.99
Brightness	\times	10%	75.0	55.4	52.0	56.7	52.0	18.14
Gaussian	\times	10%	75.2	55.7	52.1	55.4	52.2	17.29

5.2. Evaluation in Digital World

Variable-size Trigger. Generally, the data poisoning rate Poi is the key parameter to a backdoor attack. Thus we use different Poi (50%, 20%, 10%, 5%, 2%, and 1%) to evaluate the performance of our method. Besides, to prove the better performance of our attack method, we select some physical backdoor attacks on image classification as a comparison model and employ them on object detection. We find that these physical backdoor attacks only focus on single and physical transformations applying to the poisoned image, such as distance, rotation, brightness, and gaussian noise. Tab. 1 shows ASR and clean accuracy of the different backdoor object detectors based on the variable-size trigger.

In Tab. 1, AP_a has no data because there is no person in the $D_{val,a}$, we exclude the presence of “person” in $D_{val,a}$. We prepare different Poi to evaluate the impact of Poi on the performance of the backdoor object detector. The higher Poi , the lower AP_b of the backdoor object detector. It indicates that with more poisoned images, the detector learns the backdoor features better which can make more B -boxes disappear. Because the locations of the trigger do not obscure the clean features and instead provide backdoor features, the AP_b stays around 75% and the AP_{a+b} is proportional to the poisoning rate.

Compared with BadDets [19] on the same poison rate, our backdoor detector causes a 0.8% drop in AP_b for the clean images. On the other hand, ASR of our backdoor detector on backdoor images is 85.86% which is actually 65.51% higher than ASR of the BadDets on the backdoor images. Dangerous Cloaking

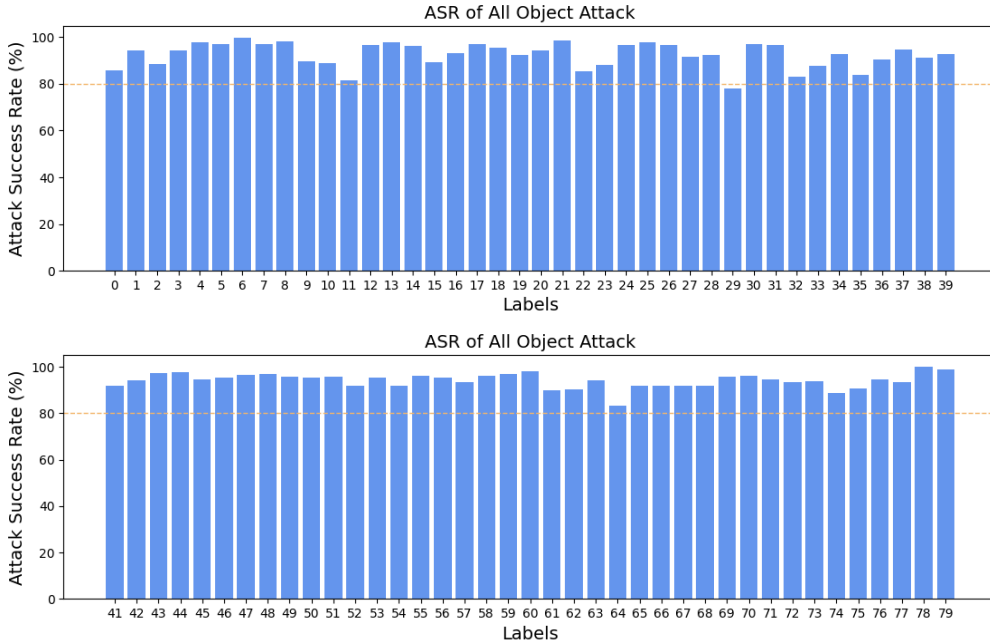


Figure 6: The performance of All Object Attack on different class. The horizontal axis is the class number in COCO dataset. The vertical axis is the ASR of backdoor object detector(%).

[22] had no experiments on the COCO dataset but on their video and trained on extra training data, so we can not get the training data to evaluate *ASR* and the clean accuracy of their detector on the COCO dataset. Besides, the ASR of physical transformations also proves that physical backdoor attacks on image classification have no effort on object detection. Because the space of the single transformation is a small subspace of big physical space, the backdoor detector learns the feature of variable triggers inadequately.

In addition, we also find if the target class is confined to only one class by the attacker during training, the trigger x_t can not cause the misdetection of the other classes in the inference phase. Therefore, to verify that the great performance of the backdoor detector depends on our attack method not the target class, we break the restriction on the specified target class and propose the All Object Attack which poisons the object containing all classes of the dataset. As Fig. 6 shows, the ASR of All Object Attack on different classes is above 80% mostly. Clearly, the performance of our attack method yields little difference across target classes.

Malicious Adversarial Training. To show the performance of robust backdoor object detector *RD*, we apply the physical noises on two-dimensional images, such as random noise, motion blur, rain and light, which create an environment that meets our needs. It validates that malicious adversarial training improves our backdoor object detector’s robustness on physical perturbations. In addition, we think BadDets has the undesirable performance without the interference of the physical noise, the performance of BadDets will be worse facing the physical test. Hence we replace BadDets with our backdoor object detector (BOD) $\mathbb{F}_{\hat{\theta}}$ as

Table 2: The results (%) of the clean detector YOLOv5s and different backdoor object detector with trigger + random noise $N(0, \sigma^2)$

Object	Metric	AP_{a+b}					mAP_{a+b}					ASR				
		0.0	0.1	0.2	0.3	0.4	0.0	0.1	0.2	0.3	0.4	0.0	0.1	0.2	0.3	0.4
YOLOv5s	σ^2	60.5	58.6	58.7	60.6	62.6	52.8	52.5	52.5	52.4	52.4	-	-	-	-	-
BOD		22.6	21.9	25.2	33.2	40.2	52.1	51.8	51.7	51.7	51.5	83.07	83.78	78.88	58.70	48.07
Distance		56.3	57.1	57.4	59.0	61.1	52.0	51.8	51.8	51.8	51.7	21.38	25.61	25.29	24.51	23.24
Rotation	$N(0, \sigma^2)$	56.2	53.0	53.4	55.6	57.1	52.1	51.9	51.9	51.8	51.7	23.99	24.24	23.82	23.27	22.34
Brightness		56.7	55.1	55.7	57.6	59.7	52.0	51.8	51.8	51.7	51.7	18.14	20.33	19.70	19.17	17.01
Gaussian		55.4	54.0	54.8	56.7	57.9	52.2	52.0	52.1	52.1	52.1	17.29	23.45	24.15	23.84	23.45
<i>RD</i> (Ours)		25.1	21.6	22.4	27.1	37.8	52.8	51.9	51.9	51.8	51.6	81.03	85.53	85.05	82.14	66.75

Table 3: The results (%) of the clean detector YOLOv5s and different backdoor object detector with trigger + motion blurring

Object	Metric	AP_{a+b}					mAP_{a+b}					ASR				
		5	10	20	50	75	5	10	20	50	75	5	10	20	50	75
YOLOv5s	degree	60.5	60.4	60.8	66.2	68.7	52.8	52.7	52.7	52.7	53.0	-	-	-	-	-
BOD		37.7	38.4	38.8	65.0	68.0	52.4	52.4	52.3	52.3	52.2	78.55	78.32	76.28	21.19	10.38
Distance	Motion Blurring	56.2	55.7	56.1	61.9	64.1	52.1	52.1	52.1	52.0	52.1	26.92	27.11	26.50	19.11	17.74
Rotation		54.3	53.9	54.6	60.6	62.1	52.2	52.1	52.1	52.1	52.2	23.45	23.35	23.28	18.34	18.02
Brightness		56.5	56.4	56.8	62.3	63.9	52.0	51.9	52.0	51.9	52.1	23.97	23.59	23.38	18.10	17.85
Gaussian		55.3	55.0	56.1	61.7	63.3	52.2	52.1	52.2	52.2	52.4	22.68	22.70	21.96	17.60	17.30
<i>RD</i> (Ours)		26.2	27.0	25.5	60.2	66.2	52.2	52.2	52.2	52.2	52.1	82.69	82.64	81.46	34.06	15.37

the normal backdoor object detector to compare with *RD* in this subsection.

When the trigger without physical noise present on object, $\mathbb{F}_{\hat{\theta}}$ and *RD* make the *B*-box of the attacked object disappear. However, when the trigger is perturbed by physical noise, only the detection of *RD* disappear successfully.

In Tab. 2, under different variances of random perturbation, in comparison with YOLOv5s and backdoor object detector, robust backdoor object detection *RD* suffering from a little performance degradation on $D_{val, a+b}$ with variance $\sigma^2 = 0$ while gaining the robustness on bigger σ^2 . When σ^2 changes from 0.0 to 0.1, *ASR* of both detectors increases a little accidentally. The reason is that the existence of random noise also influences the detector and makes a few numbers of *B*-box disappear. As σ^2 gradually increases, *ASR* of the backdoor object detector decreases greatly. In contrast, *ASR* of robust backdoor object detection decreases

Table 4: The results (%) of the clean detector YOLOv5s and different backdoor object detector with trigger + rainy

Object detector	Metric value	AP_{a+b}					mAP_{a+b}					ASR				
		50	100	150	200	250	50	100	150	200	250	50	100	150	200	250
YOLOv5s		60.7	60.8	61.0	61.1	61.1	52.7	52.7	52.8	52.7	52.8	-	-	-	-	-
BOD		41.6	42.8	43.3	44.4	45.0	52.4	52.4	52.4	52.4	52.4	74.08	72.38	71.32	70.23	69.38
Distance		56.8	57.0	57.1	57.1	57.0	52.1	52.1	52.1	52.1	52.1	26.28	25.83	25.76	25.92	25.45
Rotation	Rainy	54.5	54.5	55.3	55.3	55.5	52.1	52.1	52.1	52.1	52.1	24.06	22.58	21.97	21.91	20.27
Brightness		57.3	57.3	57.4	57.7	57.8	52.0	52.0	52.0	52.0	52.0	23.31	22.90	23.09	22.93	22.49
Gaussian		55.5	55.7	55.6	55.9	56.1	52.2	52.3	52.3	52.3	52.3	23.84	23.63	23.66	23.79	23.38
<i>RD</i> (Ours)		31.7	32.8	33.4	34.2	34.5	52.1	52.2	52.1	52.2	52.2	78.74	78.15	77.15	76.74	76.25

Table 5: The results (%) of the clean detector YOLOv5s and different backdoor object detector with trigger + light. S denote as the change times of original saturation.

Object detector	Metric S	AP_{a+b}					mAP_{a+b}					ASR				
		0.2	0.4	0.6	0.8	1.2	0.2	0.4	0.6	0.8	1.2	0.2	0.4	0.6	0.8	1.2
YOLOv5s		58.0	60.1	60.6	60.6	57.5	49.2	52.1	52.7	52.8	40.4	-	-	-	-	-
BOD		49.5	43.3	41.3	40.5	51.6	48.2	51.4	52.2	52.2	40.4	25.37	35.74	38.58	55.41	15.98
Distance		53.9	56.4	57.0	57.1	53.9	47.6	51.3	51.9	52.1	40.9	13.05	15.05	15.80	16.14	13.33
Rotation	Light	51.2	53.7	54.3	54.5	51.5	47.8	51.3	52.0	52.1	40.7	11.37	12.45	12.50	12.70	12.38
Brightness		56.5	56.5	57.1	57.2	53.8	47.7	51.1	51.8	52.0	40.9	14.98	15.45	15.37	15.28	14.24
Gaussian		53.3	55.4	55.5	55.6	52.9	47.7	51.3	51.9	52.2	40.9	10.59	11.47	12.41	12.44	11.51
<i>RD</i> (Ours)		47.8	39.3	34.9	31.7	40.5	48.0	51.6	52.3	52.3	47.2	45.08	58.99	64.71	67.81	37.95

slowly and keeps a high value when subjected to random perturbations. But the random perturbations with $\sigma = 0.4$ exceed the trigger and lead to attack failure.

As Tab. 3 shows, the degree is used to measure the fuzziness of the motion blurring on images. Even if the trigger is disturbed by a bigger degree of the motion blurring, ASR of RD still above 80%, higher than ASR of $\mathbb{F}_{\hat{\theta}}$. When the degree increases to 50, motion blurring cut the ASR down to 34.03%. The reason is that the pixel features of the trigger are severely corrupted so that both detectors can not extract the feature of the trigger making the backdoor attack fail.

In Tab. 4, the value is presented as the number of raindrops. The table shows that RD can resist the disturbance of rain. ASR of RD also is higher than the backdoor object detector in rain. But there is not much difference between them. The mAP_{a+b} shows that the rain can not disturb the backdoor attack

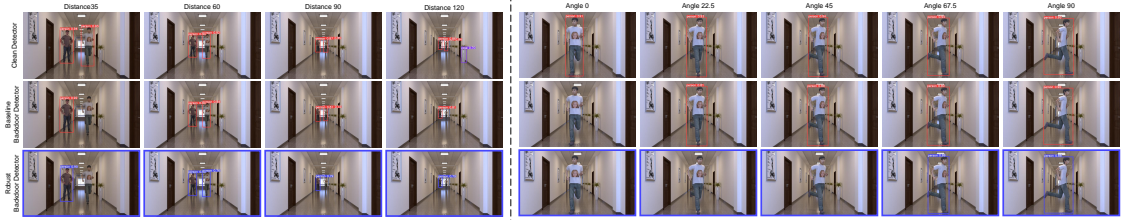


Figure 7: Visualization of different detectors with different distances and angles. The column represents the different values of distance (dm) on the left side of the black dotted line, and the right side's column represents the different angles ($^{\circ}$) of the person facing the camera.

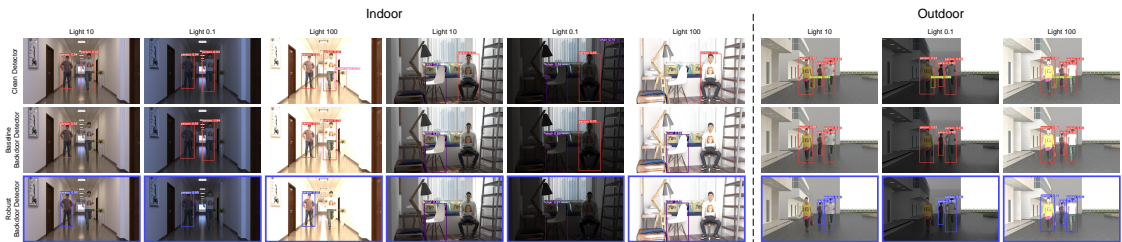


Figure 8: Visualization of different detectors with different light sources. The column represents the different light intensity values of sources. The left side of the black dotted line is the indoor environment, and the right side is the outdoor environment.

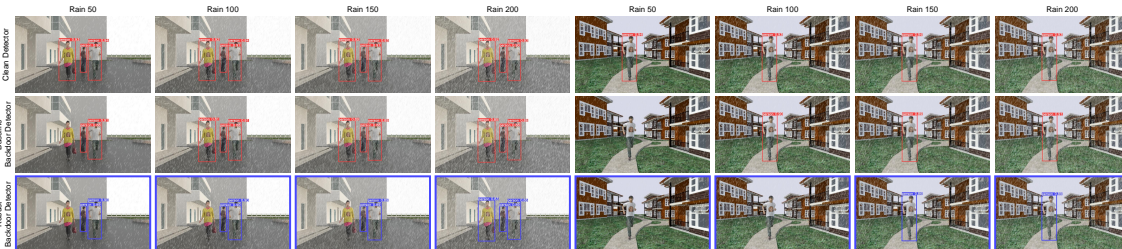


Figure 9: Visualization of the different detectors with different rains. The column represents the different number of raindrops per square meter.

severely because the raindrops are relatively evenly distributed in the image and do not completely obscure the important features of the trigger.

In Tab. 5, we change the saturation of all images to simulate the change of light in the digital world. From the table shown, our *RD* perform better than baseline backdoor object detectors. The decrease of AP_{a+b} illustrate the saturation not only influence the backdoor task but also the clean task. When the saturation change to 1.2 times, the light enable to destroy the pixels that have already high saturation. Therefore the performance on high saturation is too poor to misdetect the attacked object.

We test the trigger on the COCO and visualize the detection shown in Fig. 5. The person without a trigger has the same detection when detected by both backdoor object detectors. It indicates that our trigger does not affect the prediction of the clean rain object. Instead, when the trigger is added to the attacked object, the trigger is responded to the backdoor object detector by disappearing the *B*-box of the attacked

Table 6: The attack performance of detectors in different virtual scenes.

Object detector	Distance			Rotation			Brightness			Brightness		
	20	70	120	0	30	60	0.1	10	100	50	100	150
YOLOv5s	0.00	0.00	8.33	0.00	0.00	0.00	4.16	0.00	0.00	0.00	0.00	0.00
Baseline	91.66	16.66	12.50	91.66	20.83	4.16	4.16	79.16	8.33	100.00	75.00	33.33
<i>RD</i> (Ours)	100.0	83.33	79.16	100.0	75.00	29.16	87.5	95.83	58.33	100.00	83.33	66.66

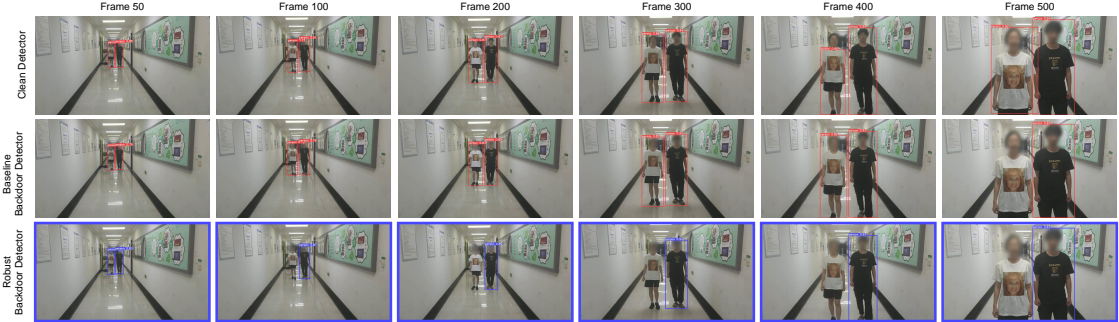


Figure 10: Visualization of different detectors indoors. The figure shows the effectiveness of trigger T-shirt for person to evading the backdoor object detector indoors. Each row corresponds to various detectors while each column shows an individual frame in a video.

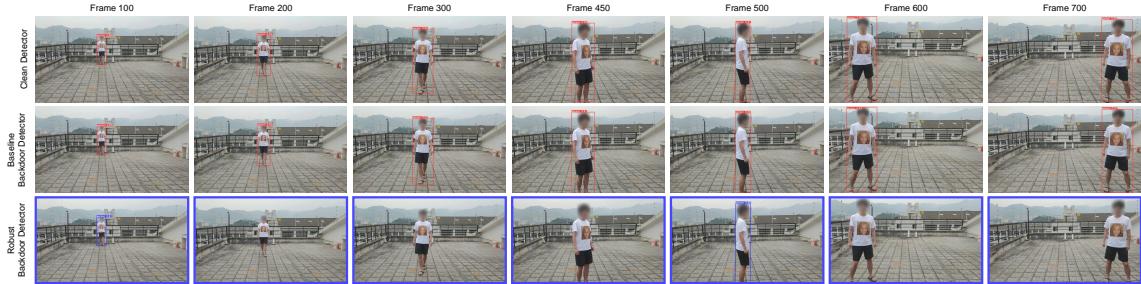


Figure 11: Visualization of different detectors outdoors in sunny. The figure shows the effectiveness of trigger T-shirt for person to evading the backdoor object detector outdoors.

object. We add the physical noises to the images to evaluate the different detector trained by our RBA and a normal backdoor attack. We can observe that RBK has superior attacking capability while another attack can not depress the impact of the physical factors.

5.3. Evaluation in Virtual World

Following the setup in the digital world, we select different human models as the attacked object. Without loss of generality, we choose background as the poisoned label to fool detectors in our experiment. We use four different physical factors illustrated in Fig. 4 for evaluating the efficacy of the proposed RBK.

As shown in Fig. 7, we find that the higher rotation angle and distance can make our *RD* less effective. This can be attributed to the fact that the trigger is badly captured when the pixel information of the trigger

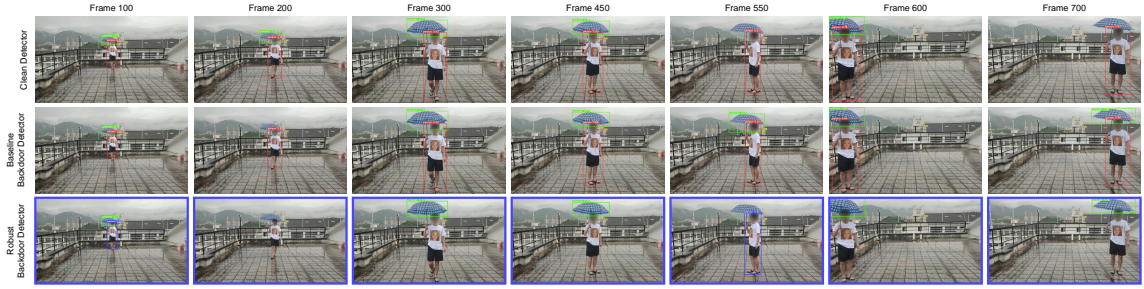


Figure 12: Visualization of different detectors outdoors in rain. The figure shows the effectiveness of trigger T-shirt for person to evading the backdoor object detector in rain.

Table 7: The attack performance of detectors in different real scene.

Object detector	Indoor, Sunny			Outdoor, Sunny			Outdoor, Rainy		
	Video 1	Video 2	Video 3	Video 4	Video 5	Video 6	Video 7	Video 8	Video 9
YOLOv5s	5.02	1.35	2.68	3.70	2.56	0.59	2.58	0.63	1.27
Baseline	17.63	21.84	6.96	5.25	3.47	4.86	5.02	4.41	6.11
<i>RD</i> (Ours)	67.65	69.01	40.87	62.56	61.42	35.94	61.42	78.10	66.88

is lost. We adjust the light sources to simulate the different daily times indoors and outdoors in Fig. 8. As the light condition is brighter or darker, the person with the trigger also is undetected, which indicates that the simple light change is invalid for our *RD*. Fig. 9 shows that *RD* is robust to the rain. But when the number of raindrops comes to 150 and 200, the person with the trigger is partially occluded in the sight of the detector. Tab. 6 shows that the percentage of misdetection in the virtual world. We find that our attack method outperforms the previous work on backdoor object detector. In addition, the robustness of *RD* is poor in the brighter scenes. This can be attributed to the fact that too bright environment can damage the most feature of the trigger captured badly by backdoor object detector.

We find the performance of *RD* is better than normal backdoor object detection when exposed to most physical factors in the virtual world.

5.4. Evaluation in Real World

To evaluate the effectiveness of *RD* in the real physical world, we make a T-shirt with a *Face* trigger. During the physical experiment in the real world, we use a smartphone to record several pieces of video. Because the distance between the viewing point and attacked object can influence the success rate of a backdoor attack, we examine the impact of varying distances and angles in Figures 10, 11 and 12. The Tab. 7 also shows the percentage of the misdetection in video set. It shows that the smaller the distance, the higher the success rate. Generally, a higher success rate is achieved at narrow angles than wide angles, and at short distances than long distances. Since *RD* has the stronger capability to recognize triggers as

Table 8: The performance (%) of the trained backdoor object detector by different trigger sizes on the three COCO datasets. The trigger transparency is fixed to 1.

Method	Trigger Size	AP_b	mAP_b	mAP_a	AP_{a+b}	mAP_{a+b}	ASR
Invariable Trigger size	20×20	75.4	55.2	51.8	50.0	51.8	50.42
	40×40	74.1	55.1	51.9	52.7	51.9	39.17
	60×60	74.7	55.2	52.0	53.5	52.0	40.49
	80×80	75.2	55.3	52.2	55.6	52.2	39.48
	100×100	75.1	55.3	52.2	57.7	52.2	33.34
Variable Trigger size	—	74.1	54.9	52.5	22.6	52.2	85.86

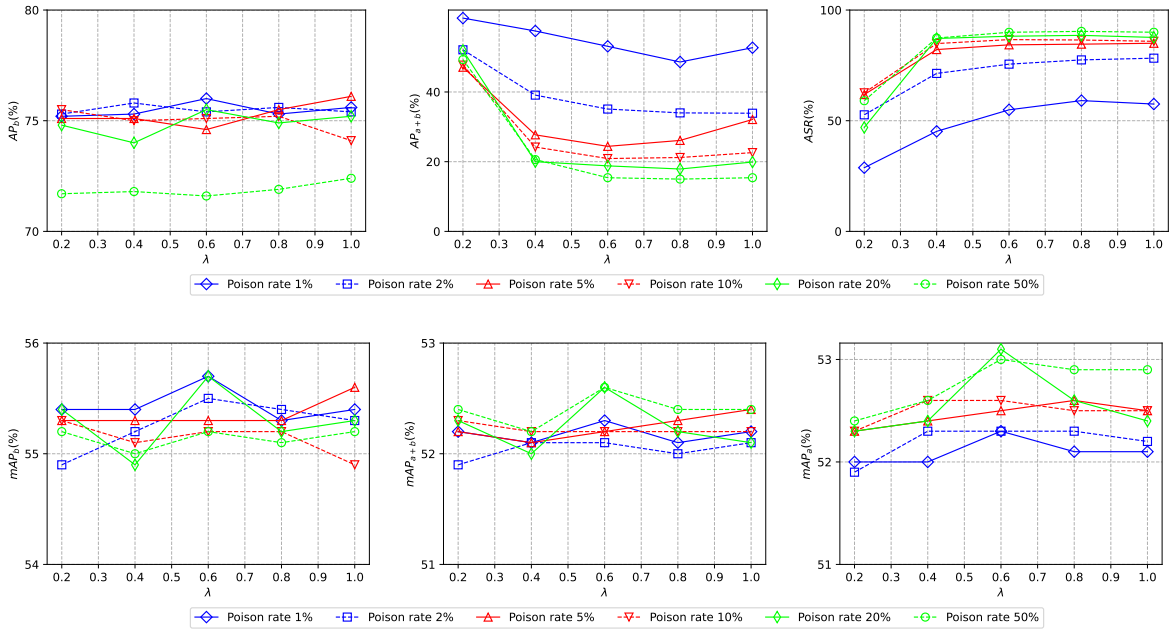


Figure 13: The clean accuracy and ASR of the backdoor object detector with different poison rate Poi and transparency λ .

the distance comes smaller, it is easier for RD to hide the B -box. The performance of the baseline backdoor object detector is weaker than RD when facing numerous physical factors. Once the proposed RBK method has been exploited during the training of object detector, the performance of backdoor object detector are greatly improved.

Table 9: Training YOLOv5s with different backbone, clean accuracy (%) and *ASR* (%) of backdoor attacks with 10% poisoning rate and transparency of 1.

Object detector	Backbone	Poison Rate	AP_b	mAP_b	mAP_a	AP_{a+b}	mAP_{a+b}	<i>ASR</i>
YOLOv5s	DarkNet	10%	74.1	54.9	52.5	22.6	52.2	85.86
		0%	76.4	56.7	52.7	60.5	52.8	–
	VGG11 [49]	10%	76.3	55.5	52.5	23.0	52.5	84.82
		0%	76.2	55.5	51.8	60.1	51.9	–
	ResNet50 [2]	10%	76.9	59.4	57.2	18.5	56.7	87.15
		0%	76.6	58.7	55.2	61.3	55.2	–
	MobileNetV3 [50]	10%	68.0	47.1	44.9	22.4	44.7	84.17
		0%	68.7	46.8	44.0	56.8	44.2	–
	DenseNet121 [51]	10%	75.7	57.0	54.7	22.2	54.3	86.68
		0%	76.2	56.5	52.9	54.3	52.6	–

5.5. Ablation Analysis

In this section, we do a series of experiments on the COCO dataset. We study the impact of different parameters on *ASR* and the clean accuracy of the backdoor object detector: (1) backdoor trigger size, (2) transparency λ , (3) backbone, and (4) loss function L_v . In the next experiments, we set the trigger as *Face*. Also, for the study on each parameter of RBK, we fix the other parameters.

Ablation Study on Trigger Size. Here, we show how the invariable-size triggers impact the clean accuracy and *ASR* of the backdoor object detector. Different trigger sizes consist of 20×20 to 120×120 at 20 intervals. As shown in Tab. 8, the 20×20 size of the trigger achieve 75.4% AP_b of the “person” class and 50.42% *ASR*. With the size of the trigger from 20×20 to 120×120 , *ASR* of the backdoor attack decreases from 50.42% down to 31.15%. In contrast, the variable-size trigger reaches 85.86% *ASR*. Despite the changes of AP_a and mAP_b , the clean accuracy of the variable-size trigger is not much different from the others. This ablation study shows that the performance improvement using variable-size triggers is significant.

Ablation Study on λ . Whether the triggers are human-perceived during training is a critical issue. Transparency enhances the invisibility of the triggers when added to the training dataset. Fig. 13 shows the variation of clean accuracy and *ASR* on different transparency λ , and poison rate *Poi*. In general, the larger λ is, the easier for a backdoor to attack the detector, and the bigger *ASR* of the backdoor attack is. Of interest are the results for the backdoor object detector, *ASR* decrease a little from 85% to 81% during $Poi \geq 5\%$ and decrease rapidly to 53% with only a slight drop on *Poi*. It can be seen that the change

Table 10: Clean accuracy (%) and *ASR* (%) of Robust Backdoor training using different loss function L_y and L_v .

Method	AP_b	mAP_b	mAP_a	AP_{a+b}	mAP_{a+b}	<i>ASR</i>
YOLO L_v, L_y	74.1	54.9	52.5	22.6	52.2	85.86
YOLO w/o L_v	75.8	55.8	52.7	31.9	52.7	80.12
YOLO w/o L_y	73.7	54.7	52.4	20.1	51.9	87.15

of both parameters will have little influence on the clean accuracy except when the poison rate is 50%, the level of poison rate makes serious damage to AP_b which reduces by 3%. Fig. 13 also shows mAP_b and mAP_a illustrating that poison rates and transparency have little effect on the average accuracy of other classes, which fluctuate between 0.4%. The sensitivity of its mAP_{a+b} to the poison rate and transparency is even weaker in the dataset $D_{val,a}$. Given the above, Poi and λ also make an impact on *ASR* and have little influence on clean accuracy.

Ablation Study on Backbone. The structure of an object detector is usually composed of different components. The extracted feature of the trigger depends on the backbone which is one component of the object detector. So we show whether the different backbones affect the performance of the backdoor attack. We compare different backbone VGG11 [49], Resnet [2], MobileNetV3 [50], DenseNet121 [51] for ablation experiments in Tab. 9.

The clean accuracy and *ASR* of five backbones are ordered from smallest to largest as MobileNetV3, VGG11, DarkNet, DenseNet121, and ResNet50. Darknet is the backbone that has small parameter quantities achieving large accuracy. It can be seen that the larger parameter quantities which strengthen the model’s ability to recognize the trigger will cause the model easier to be attacked more. ResNet50 which makes the trigger better recognized outperforms another backbone with 2.5% higher clean accuracy and 1.3% higher *ASR* than Darknet. But the training time increases to 2.15 times and the inference time increases to 1.38 times at the cost of increased accuracy.

Ablation Study on L_v and L_y . The strong perturbations may cause interference on the clean detector and affect the backdoor attack, so we introduce L_v and L_y to improve the clean accuracy and *ASR*. Tab. 10 shows that the performance of the backdoor object detector without L_v or L_y decreases on both clean and poisoned images. Without L_v , *ASR* of robust backdoor object detector *RD* decrease by 5.75%, and the AP_b increase by 1.7% respectively. The loss function L_v makes the stronger perturbation to be learned by the detector so that *RD* has a greater *ASR*. Without L_y , the clean accuracy of the backdoor object detector decrease by 0.4%, and *ASR* increase by 1.29%.

Table 11: The percentage of B -box (%) detected as the person class by clean detector, Backdoor object detector, and Robust backdoor object detector with different datasets. Low represents $Score_B \in [0, 0.1]$. Middle represents $Score_B \in (0.1, 0.5]$. High represents $Score_B \in (0.5, 1.0]$.

Object detector	Dataset	Clean			Backdoor			Backdoor with Δ_{phy}			
		$Score_B$	Low	Middle	High	Low	Middle	High	Low	Middle	High
YOLOv5s	B -box Quantity		99.55	0.31	0.14	99.43	0.46	0.11	99.49	0.41	0.10
Baseline			99.35	0.49	0.16	99.68	0.30	0.02	99.53	0.38	0.09
RD (Ours)			99.36	0.48	0.16	99.68	0.29	0.03	99.64	0.32	0.04

6. Discussion

In this section, we discuss the possible reasons why our attack method is effective. When disturbed by the physical factor occurring in the real world, the trigger will lose connection with backdoor-related neurons of backdoor object detection. The blocked feature of the trigger will be damaged which can respond to the backdoor-related neurons. Therefore, the original B -box with the higher $Score_B$ will be retained by NMS. We can observe that the detection of the object with the trigger is y , not \hat{y} .

Therefore, the adversary hopes the backdoor object detector learns about more different sizes of triggers with more physical noises. By learning the feature of this trigger hard to be extracted by the backbone, RD reduces the sensitivity to physical factors. The backdoor-related neurons will be activated which affect the generation of the original B -box by lowering the $Score_B$ of the original B -box seen by Tab. 11. Therefore, the number of B -box closed to y will be limited and the number of B -box closed to \hat{y} will be retained by NMS. As a result, the object with the trigger disturbed by physical factors will be detected as \hat{y} .

To summarize, our RBK increases the diversity of backdoor triggers, making the backdoor object detector strengthen the association between these triggers and poisoned labels. It widens the boundary where an object with a trigger will be detected as a poisoned label. Even if the physical factors are added to the trigger, the object with the trigger will be near the boundary but still stay in the boundary of RD .

7. Conclusion

This paper introduces the robust backdoor attack on object detectors. We report the disadvantages of existing backdoor attacks on object detection, which lack robustness to physical factors and improve the performance only in the digital world. We first propose a variable-size trigger that can fit different sizes of attacked objects to reflect the distance between the viewing point and attacked object in the real world. In addition, to enhance the robustness of the backdoor object detector on physical factors, we propose malicious adversarial training to adapt the detector to most physical factors, for which we generate the strongest physical noise example hard to be detected by the detector in the pixel space. Experiments demonstrate that our method performs effectively in the digital, virtual, and real worlds. We also demonstrate that our

method ensures the robustness of the backdoor object detector on physical noises and vertical rotation of objects. In summary, security-sensitive domains need to focus on the security threat of robust backdoor attacks. Backdoor attacks on object detection robust to the real world can be harmful to humans.

Acknowledgment

This work was supported by the National Natural Science Foundation of China (92167203), Zhejiang Provincial Natural Science Foundation of China (LZ22F020007), National Natural Science Foundation of China (62202457), China Postdoctoral Science Foundation (2022M713253), and Science and Technology Innovation Foundation for Graduate Students of Zhejiang University of Science and Technology (F464108M05).

References

- [1] X. Cao, C. Li, J. Feng, L. Jiao, Semi-supervised feature learning for disjoint hyperspectral imagery classification, *Neurocomputing*.
- [2] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, IEEE Computer Society, 2016, pp. 770–778. doi:10.1109/CVPR.2016.90.
- [3] S. Huang, T. Wang, H. Xiong, B. Wen, J. Huan, D. Dou, Temporal output discrepancy for loss estimation-based active learning, CoRR abs/2212.10613. arXiv:2212.10613, doi:10.48550/arXiv.2212.10613.
- [4] H. Zhang, J. Wang, Towards adversarially robust object detection, in: 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019, IEEE, 2019, pp. 421–430. doi:10.1109/ICCV.2019.00051.
- [5] J. Yang, S. Liu, Z. Li, X. Li, J. Sun, Real-time object detection for streaming perception, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, IEEE, 2022, pp. 5375–5385. doi:10.1109/CVPR52688.2022.00531.
- [6] K. J. Joseph, S. H. Khan, F. S. Khan, V. N. Balasubramanian, Towards open world object detection, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE, 2021, pp. 5830–5840. doi:10.1109/CVPR46437.2021.00577.
- [7] T. Zhou, W. Wang, S. Liu, Y. Yang, L. V. Gool, Differentiable multi-granularity human representation learning for instance-aware human semantic parsing, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE, 2021, pp. 1622–1631. doi:10.1109/CVPR46437.2021.00167.
- [8] S. Qiu, S. Anwar, N. Barnes, Semantic segmentation for real point cloud scenes via bilateral augmentation and adaptive fusion, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE, 2021, pp. 1757–1767. doi:10.1109/CVPR46437.2021.00180.
- [9] M. Fan, S. Lai, J. Huang, X. Wei, Z. Chai, J. Luo, X. Wei, Rethinking bisenet for real-time semantic segmentation, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, Computer Vision Foundation / IEEE, 2021, pp. 9716–9725. doi:10.1109/CVPR46437.2021.00959.
- [10] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. S. Bernstein, A. C. Berg, L. Fei-Fei, Imagenet large scale visual recognition challenge, *Int. J. Comput. Vis.* 115 (3) (2015) 211–252. doi:10.1007/s11263-015-0816-y.

- [11] J. H. Metzen, T. Genewein, V. Fischer, B. Bischoff, On detecting adversarial perturbations, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings, OpenReview.net, 2017.
- [12] W. Xu, D. Evans, Y. Qi, Feature squeezing: Detecting adversarial examples in deep neural networks, in: 25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018, The Internet Society, 2018.
- [13] Z. You, J. Ye, K. Li, Z. Xu, P. Wang, Adversarial noise layer: Regularize neural network by adding noise, in: 2019 IEEE International Conference on Image Processing, ICIP 2019, Taipei, Taiwan, September 22-25, 2019, IEEE, 2019, pp. 909–913. doi:10.1109/ICIP.2019.8803055.
- [14] T. Gu, B. Dolan-Gavitt, S. Garg, Badnets: Identifying vulnerabilities in the machine learning model supply chain, CoRR abs/1708.06733. arXiv:1708.06733.
- [15] Z. Wang, J. Zhai, S. Ma, Bppattack: Stealthy and efficient trojan attacks against deep neural networks via image quantization and contrastive adversarial learning, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, IEEE, 2022, pp. 15054–15063. doi:10.1109/CVPR52688.2022.01465.
- [16] J. Zhang, D. Chen, Q. Huang, J. Liao, W. Zhang, H. Feng, G. Hua, N. Yu, Poison ink: Robust and invisible backdoor attack, IEEE Trans. Image Process. 31 (2022) 5691–5705. doi:10.1109/TIP.2022.3201472.
- [17] X. Chu, A. Zheng, X. Zhang, J. Sun, Detection in crowded scenes: One proposal, multiple predictions, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, Computer Vision Foundation / IEEE, 2020, pp. 12211–12220. doi:10.1109/CVPR42600.2020.01223.
- [18] S. M. Ahmed, A. R. Lejbølle, R. Panda, A. K. Roy-Chowdhury, Camera on-boarding for person re-identification using hypothesis transfer learning, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, Computer Vision Foundation / IEEE, 2020, pp. 12141–12150. doi:10.1109/CVPR42600.2020.01216.
- [19] S. Chan, Y. Dong, J. Zhu, X. Zhang, J. Zhou, Baddet: Backdoor attacks on object detection, in: L. Karlinsky, T. Michaeli, K. Nishino (Eds.), Computer Vision - ECCV 2022 Workshops - Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part I, Vol. 13801 of Lecture Notes in Computer Science, Springer, 2022, pp. 396–412. doi:10.1007/978-3-031-25056-9_26.
- [20] H. Ma, Y. Li, Y. Gao, Z. Zhang, A. Abuadba, A. Fu, S. F. Al-Sarawi, S. Nepal, D. Abbott, MACAB: model-agnostic clean-annotation backdoor to object detection with natural trigger in real-world, CoRR abs/2209.02339. arXiv:2209.02339, doi:10.48550/arXiv.2209.02339.
- [21] T. Wu, T. Wang, V. Schwag, S. Mahloujifar, P. Mittal, Just rotate it: Deploying backdoor attacks via rotation transformation, in: A. Demontis, X. Chen, F. Tramèr (Eds.), Proceedings of the 15th ACM Workshop on Artificial Intelligence and Security, AISec 2022, Los Angeles, CA, USA, 11 November 2022, ACM, 2022, pp. 91–102. doi:10.1145/3560830.3563730.
- [22] H. Ma, Y. Li, Y. Gao, A. Abuadba, Z. Zhang, A. Fu, H. Kim, S. F. Al-Sarawi, S. Nepal, D. Abbott, Dangerous cloaking: Natural trigger based backdoor attacks on object detectors in the physical world, CoRR abs/2201.08619. arXiv:2201.08619.
- [23] A. Shrivastava, A. Gupta, R. B. Girshick, Training region-based object detectors with online hard example mining, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, IEEE Computer Society, 2016, pp. 761–769. doi:10.1109/CVPR.2016.89.
- [24] D. Wu, Y. Wang, Adversarial neuron pruning purifies backdoored deep models, in: M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang, J. W. Vaughan (Eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, 2021, pp. 16913–16925.
- [25] Y. Liu, S. Ma, Y. Aafer, W. Lee, J. Zhai, W. Wang, X. Zhang, Trojaning attack on neural networks, in: 25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-21, 2018,

The Internet Society, 2018.

- [26] X. Chen, C. Liu, B. Li, K. Lu, D. Song, Targeted backdoor attacks on deep learning systems using data poisoning, CoRR abs/1712.05526. [arXiv:1712.05526](https://arxiv.org/abs/1712.05526).
- [27] H. Zhong, C. Liao, A. C. Squicciarini, S. Zhu, D. J. Miller, Backdoor embedding in convolutional neural network models via invisible perturbation, in: V. Roussev, B. Thuraisingham, B. Carminati, M. Kantarcioglu (Eds.), CODASPY '20: Tenth ACM Conference on Data and Application Security and Privacy, New Orleans, LA, USA, March 16-18, 2020, ACM, 2020, pp. 97–108. doi:10.1145/3374664.3375751.
- [28] Q. Zhang, Y. Ding, Y. Tian, J. Guo, M. Yuan, Y. Jiang, Advdoor: adversarial backdoor attack of deep learning system, in: C. Cadar, X. Zhang (Eds.), ISSTA '21: 30th ACM SIGSOFT International Symposium on Software Testing and Analysis, Virtual Event, Denmark, July 11-17, 2021, ACM, 2021, pp. 127–138. doi:10.1145/3460319.3464809.
- [29] S. Moosavi-Dezfooli, A. Fawzi, P. Frossard, Deepfool: a simple and accurate method to fool deep neural networks, CoRR abs/1511.04599. [arXiv:1511.04599](https://arxiv.org/abs/1511.04599).
- [30] M. Barni, K. Kallas, B. Tondi, A new backdoor attack in CNNs by training set corruption without label poisoning, in: 2019 IEEE International Conference on Image Processing, ICIP 2019, Taipei, Taiwan, September 22-25, 2019, IEEE, 2019, pp. 101–105. doi:10.1109/ICIP.2019.8802997.
- [31] Y. Liu, X. Ma, J. Bailey, F. Lu, Reflection backdoor: A natural backdoor attack on deep neural networks, in: A. Vedaldi, H. Bischof, T. Brox, J. Frahm (Eds.), Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part X, Vol. 12355 of Lecture Notes in Computer Science, Springer, 2020, pp. 182–199. doi:10.1007/978-3-030-58607-2_11.
- [32] R. Ning, J. Li, C. Xin, H. Wu, Invisible poison: A blackbox clean label backdoor attack to deep neural networks, in: 40th IEEE Conference on Computer Communications, INFOCOM 2021, Vancouver, BC, Canada, May 10-13, 2021, IEEE, 2021, pp. 1–10. doi:10.1109/INFOCOM42981.2021.9488902.
- [33] J. Dumford, W. J. Scheirer, Backdooring convolutional neural networks via targeted weight perturbations, in: 2020 IEEE International Joint Conference on Biometrics, IJCB 2020, Houston, TX, USA, September 28 - October 1, 2020, IEEE, 2020, pp. 1–9. doi:10.1109/IJCB48548.2020.9304875.
- [34] Y. Ji, X. Zhang, T. Wang, Backdoor attacks against learning systems, in: 2017 IEEE Conference on Communications and Network Security, CNS 2017, Las Vegas, NV, USA, October 9-11, 2017, IEEE, 2017, pp. 1–9. doi:10.1109/CNS.2017.8228656.
- [35] R. Tang, M. Du, N. Liu, F. Yang, X. Hu, An embarrassingly simple approach for trojan attack in deep neural networks, in: R. Gupta, Y. Liu, J. Tang, B. A. Prakash (Eds.), KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020, ACM, 2020, pp. 218–228. doi:10.1145/3394486.3403064.
- [36] Y. Li, Y. Li, Y. Lv, Y. Jiang, S. Xia, Hidden backdoor attack against semantic segmentation models, CoRR abs/2103.04038. [arXiv:2103.04038](https://arxiv.org/abs/2103.04038).
- [37] Y. Feng, B. Ma, J. Zhang, S. Zhao, Y. Xia, D. Tao, FIBA: frequency-injection based backdoor attack in medical image analysis, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, IEEE, 2022, pp. 20844–20853. doi:10.1109/CVPR52688.2022.02021.
- [38] J. Mao, Y. Qian, J. Huang, Z. Lian, R. Tao, B. Wang, W. Wang, T. Yao, Object-free backdoor attack and defense on semantic segmentation, Computers & Security 132 (2023) 103365. doi:<https://doi.org/10.1016/j.cose.2023.103365>.
- [39] E. Bagdasaryan, V. Shmatikov, Spinning language models: Risks of propaganda-as-a-service and countermeasures, in: 43rd IEEE Symposium on Security and Privacy, SP 2022, San Francisco, CA, USA, May 22-26, 2022, IEEE, 2022, pp. 769–786. doi:10.1109/SP46214.2022.9833572.
- [40] Y. Liu, G. Shen, G. Tao, S. An, S. Ma, X. Zhang, Piccolo: Exposing complex backdoors in NLP transformer models, in: 43rd IEEE Symposium on Security and Privacy, SP 2022, San Francisco, CA, USA, May 22-26, 2022, IEEE, 2022, pp.

2025–2042. doi:10.1109/SP46214.2022.9833579.

- [41] J. Lu, H. Sibai, E. Fabry, Adversarial examples that fool detectors, CoRR abs/1712.02494. arXiv:1712.02494.
- [42] S. Thys, W. V. Ranst, T. Goedemé, Fooling automated surveillance cameras: Adversarial patches to attack person detection, in: IEEE Conference on Computer Vision and Pattern Recognition Workshops, CVPR Workshops 2019, Long Beach, CA, USA, June 16-20, 2019, Computer Vision Foundation / IEEE, 2019, pp. 49–55. doi:10.1109/CVPRW.2019.00012.
- [43] M. Sharif, S. Bhagavatula, L. Bauer, M. K. Reiter, Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition, in: E. R. Weippl, S. Katzenbeisser, C. Kruegel, A. C. Myers, S. Halevi (Eds.), Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security, Vienna, Austria, October 24-28, 2016, ACM, 2016, pp. 1528–1540. doi:10.1145/2976749.2978392.
- [44] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, D. Song, Robust physical-world attacks on deep learning visual classification, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018, Computer Vision Foundation / IEEE Computer Society, 2018, pp. 1625–1634. doi:10.1109/CVPR.2018.00175.
- [45] A. Athalye, L. Engstrom, A. Ilyas, K. Kwok, Synthesizing robust adversarial examples, in: J. G. Dy, A. Krause (Eds.), Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018, Vol. 80 of Proceedings of Machine Learning Research, PMLR, 2018, pp. 284–293.
- [46] Y. Zhao, H. Zhu, R. Liang, Q. Shen, S. Zhang, K. Chen, Seeing isn't believing: Towards more robust adversarial attack against real world object detectors, in: L. Cavallaro, J. Kinder, X. Wang, J. Katz (Eds.), Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security, CCS 2019, London, UK, November 11-15, 2019, ACM, 2019, pp. 1989–2004. doi:10.1145/3319535.3354259.
- [47] K. Xu, G. Zhang, S. Liu, Q. Fan, M. Sun, H. Chen, P. Chen, Y. Wang, X. Lin, Adversarial t-shirt! evading person detectors in a physical world, in: A. Vedaldi, H. Bischof, T. Brox, J. Frahm (Eds.), Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part V, Vol. 12350 of Lecture Notes in Computer Science, Springer, 2020, pp. 665–681. doi:10.1007/978-3-030-58558-7_39.
- [48] N. Suryanto, Y. Kim, H. Kang, H. T. Larasati, Y. Yun, T. Le, H. Yang, S. Oh, H. Kim, DTA: physical camouflage attacks using differentiable transformation network, in: IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022, IEEE, 2022, pp. 15284–15293. doi:10.1109/CVPR52688.2022.01487.
- [49] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, in: Y. Bengio, Y. LeCun (Eds.), 3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings, 2015.
- [50] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, H. Adam, Mobilenets: Efficient convolutional neural networks for mobile vision applications, CoRR abs/1704.04861. arXiv:1704.04861.
- [51] G. Huang, Z. Liu, K. Q. Weinberger, Densely connected convolutional networks, CoRR abs/1608.06993. arXiv:1608.06993.