# Fooling the Eyes of Autonomous Vehicles: Robust Physical Adversarial Examples Against Traffic Sign Recognition Systems

Wei Jia
School of Cyber Science and Engineering
Huazhong Univ. of Sci. & Tech.
jiaw@hust.edu.cn

Zhaojun Lu*
School of Cyber Science and Engineering
Huazhong Univ. of Sci. & Tech.
lzj_cse@hust.edu.cn

Haichun Zhang
Huazhong Univ. of Sci. & Tech.
homer@thesimpsons.com

Zhenglin Liu
Huazhong Uni. of Sci. & Tech.
liuzhenglin@hust.edu.cn

Jie Wang
Shenzhen Kaiyuan Internet Security Co., Ltdy
wangjie@seczone.cn

Gang Qu
University of Maryland
gangqu@umd.edu

*Abstract*—Adversarial Examples (AEs) can deceive Deep Neural Networks (DNNs) and have received a lot of attention recently. However, majority of the research on AEs is in the digital domain and the adversarial patches are static. Such research is very different from many real-world DNN applications such as Traffic Sign Recognition (TSR) systems in autonomous vehicles. In TSR systems, object detectors use DNNs to process streaming video in real time. From the view of object detectors, the traffic sign's position and quality of the video are continuously changing, rendering the digital AEs ineffective in the physical world.

In this paper, we propose a systematic pipeline to generate robust physical AEs against real-world object detectors. Robustness is achieved in three ways. First, we simulate the in-vehicle cameras by extending the distribution of image transformations with the blur transformation and the resolution transformation. Second, we design the single and multiple bounding boxes filters to improve the efficiency of the perturbation training. Third, we consider four representative attack vectors, namely Hiding Attack (HA), Appearance Attack (AA), Non-Target Attack (NTA) and Target Attack (TA). For each of them, a loss function is defined to minimize the impact of the fabrication process on the physical AEs.

We perform a comprehensive set of experiments under a variety of environmental conditions by varying the distance from $0m$ to $30m$, changing the angle from $-60°$ to $60°$, and considering illuminations in sunny and cloudy weather as well as at night. The experimental results show that the physical AEs generated from our pipeline are effective and robust when attacking the YOLO v5 based TSR system. The attacks have good transferability and can deceive other state-of-the-art object detectors. We launched HA and NTA on a brand-new 2021 model vehicle. Both attacks are successful in fooling the TSR system, which could be a lifethreatening case for autonomous vehicles. Finally, we discuss three defense mechanisms based on image preprocessing, AEs detection, and model enhancing.

## I. INTRODUCTION

Artificial Intelligence (AI) and Deep Neural Networks (DNNs) have boosted the performance of a large variety of applications, in particular computer vision tasks such as face recognition [1], image classification [2], and object detection [3]. Unfortunately, the ubiquity and diversity of AI applications have also created incentives and opportunities for attackers to attack DNNs for malicious purposes [4]. In 2014, Szegedy *et al.* [5] first introduced the concept of Adversarial Examples (AEs), which are well-crafted examples with imperceptive perturbations that deceive the image classifiers into misclassification. A lot of follow-up studies [6–10] demonstrated the vulnerabilities of DNNs under different types of AEs. One common feature of these attacks is that they are all in the digital domain instead of physical domain. Compared to the theoretical studies which have yielded prolific results on AEs and led to better understanding the principles of DNNs, the practicality of these AEs and their impact on real-world AI applications remain under-investigated.

Recently there have been reported studies on the feasibility of AEs in the physical domain, where the digital adversarial images are printed to the physical domain first, and then pictures are taken from these physical images to create AEs [11–15]. However, the physical AEs generated by this digital-physical-digital conversion becomes significantly less effective because of complicated physical conditions during the process such as the distance, angle, and illumination when the images are re-taken. Consequently, there are approaches to improve the robustness. Initial efforts have been devoted to improving the robustness of AEs. Athalye *et al.* [12] introduced the Expectation Over Transformation (EOT) method to simulate the effect of rotation, scaling, and perspective changes. Evtimov *et al.* [13] and Sitawarin *et al.* [14] enhanced this method by synthesizing the adversarial traffic signs to attack the image classifiers. Object detectors are adopted for the Traffic Sign Recognition (TSR) task to assist the safety- and security-critical autonomous driving systems [15].

Unlike image classifiers, fooling object detectors with physical AEs is much more challenging for the following

reasons. Firstly, image classifiers only process static images, but the object detectors like those in autonomous vehicles are commonly working in environments where physical features of the object such as its relative position to the object detector keep changing. Secondly, the imperfections in the fabrication process have an uncontrollable impact on the effectiveness of AEs, creating a large gap between the digital domain and the physical domain for adversarial attacks. The digital-physical-digital conversion would weaken the toxicity of the fabricated AEs. Thirdly, the AEs generation algorithms require some information of the targeted DNNs, which may not be realistic because attackers may not have control over the built-in systems of autonomous vehicles. One practical approach is to generate AEs under the white-box settings and use them to attack the black-box object detectors. To summarize, *the key to fooling object detectors, the eyes of autonomous vehicles, is to improve the robustness of physical AEs.*

In this paper, we propose a systematic pipeline to generate robust physical AEs and demonstrate its effectiveness against the object detectors used in TSR systems. To reflect the real-world scenarios, we generate physical AEs for traffic signs with a large range of physical parameters including distance, angle, and illumination. First, we extend the distribution of image transformations with the blur transformation and the resolution transformation to simulate the in-vehicle cameras. Then, Single Bounding Box (S-BBOX) and Multiple Bounding Box (M-BBOX) filters are designed to obtain the relative BBOXes before generating the physical AEs. We define four loss functions to train adversarial perturbations corresponding to the following four attack vectors: Hiding Attack (HA), Appearance Attack (AA), Non-Target Attack (NTA), and Target Attack (TA). HA hides AEs in the background so that the object detectors cannot detect them. AA makes the object detectors recognize a bizarre AE as a common category. Both NTA and TA deceive the object detector into misrecognition with imperceptible AEs. TA is more destructive since it makes the object detectors recognize an AE of one category as an object from another designated category. To improve the transferability of AEs, we use different background images in the AEs generation algorithm to avoid overfitting.

Finally, we validate the robustness of our physical AEs by driving a brand-new 2021 model vehicle toward the physical AEs to see whether the vehicle's TSR system will be fooled. Our main contributions can be summarized as follows:

*1) A systematic approach to generate robust physical AEs:* In our physical AE generation pipeline, we propose several approaches to improve the robustness of the physical AEs. We extend the distribution of image transformations to simulate the complicated environmental driving conditions. S-BBOX and M-BBOX filters are designed to obtain the BBOXes associated with the target object to train perturbation efficiently. Four loss functions are defined to generate AEs for attack vectors.

*2) A comprehensive set of experiments:* To systematically evaluate the effectiveness and robustness of the generated physical AEs, we conduct extensive outdoor experiments. More than 1,000 video clips containing more than 100,000 image frames are taken by a high-resolution camera. Real driving scenarios are simulated with varying distances from $1m$ to $30m$, angles from $-60°$ to $60°$, and illuminations for sunny, cloudy, and night.

*3) Successful attacks against YOLO v5 based object detector and TSR system in a 2021 model vehicle:* To the best of our knowledge, this is the first set of adversarial attacks against YOLO v5 based object detectors in the physical domain. We successfully launch four attack vectors, especially NTA and TA, that are life-threatening in the real world. Our physical AEs also exhibit satisfactory transferability when attacking a production-grade TSR system of a brand-new 2021 model vehicle.

## II. Background

In this section, we first introduce the advances in the field of image classification and object detection and explain why autonomous vehicles adopt object detectors for the TSR task. Then, we summarize the difficulties and limitations of the existing physical adversarial attacks against the object detectors.

### A. Object Detection

An image classifier $f(\cdot) : \mathbb{R}^{h \times w \times \varepsilon} \to \mathbb{R}^N$ recognizes an input image $x \in \mathbb{R}^{h \times w \times \varepsilon}$ to a category label $y_x \in (1, ..., N)$ where $h$, $w$ and $\varepsilon$ is the height, width and channel of $x$. The image classifier is commonly trained with supervised learning strategy, whose goal is to minimize a loss function between the output of the network $f(x)$ and the expected label $y_x$ [17]. However, image classification cannot satisfy the demands of processing dynamic video streaming containing multiple objects for autonomous vehicles [18]. Take the TSR task for example, it is vital for the vehicle sensor to recognize a traffic sign quickly and accurately from the complicated road background, then give the correct instruction to the vehicle controller for rapid response [19]. Object detection based on DNNs is adopted by autonomous vehicles to process consecutive frames of images containing multi-category objects (*e.g.*, traffic signs, vehicles, pedestrians, and cyclists) [20]. Modern DNNs based object detectors can be classified into two categories: one-stage architecture represented by YOLO series with higher detecting speed [21], and two-stage architecture represented by Faster R-CNN with higher detecting accuracy [22]. Since one-stage architecture processes BBOX regression and object classification concurrently without a region proposal stage, it is much faster than two-stage architecture so that meets the real-time requirement [23]. In recent years, many algorithms have been developed to balance the speed and the accuracy of the object detectors. Bochkovskiy *et al.* [24] proposed YOLO v4 that improved the accuracy compared with YOLO v3 [25] while maintaining the real-time object detection capability. Soon after the release of YOLO v4, a company named Ultralytics released YOLO v5's source code at Github [26]. YOLO v5 has higher mean Average Precision (mAP) and lower processing time than YOLO v4. It should be emphasized that YOLO v5 is trained with clear images, but it can be used to detect blurred images [16]. The prominent features allow YOLO v5 to detect objects under various conditions, which is well-performed in the TSR task.

As briefly introduced in Fig.1, YOLO v5 inherits some advanced techniques from YOLO v4, including the Darknet53 with Cross Stage Partial [27] (CSPDarknet53) as backbone, the Path Aggregation Network [28] (PANet) as neck, and the Spatial Pyramid Pooling [29] (SPP) block to increase the
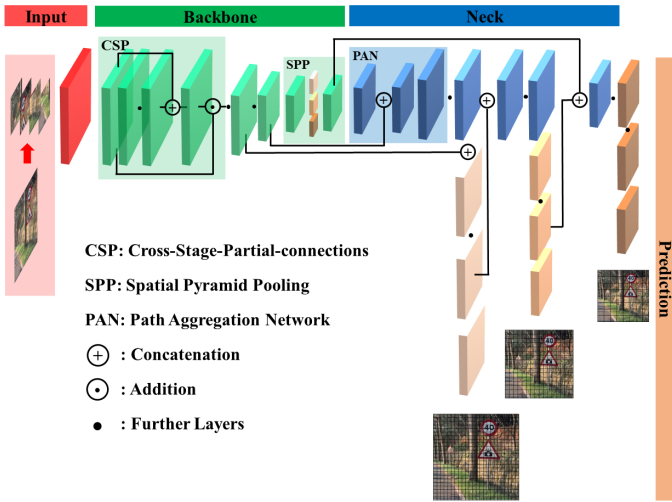
Fig. 1. Architecture of YOLO v5.



Fig. 2. The Digital-Physical-Digital cross-domain conversion of AEs.

receptive field. Besides, YOLO v5 makes some fine-tuning on the basis of YOLO v4 such as slicing input images into four small ones and concatenating them together for convolution operation. As a classical anchor-base object detector, YOLO v5 first pre-processes the anchors, then corrects the Bounding Box (BBOX) according to the off-set values, finally obtains the probability of the detected object in this BBOX [24]. Each anchor outputs three parts, the object probability refers to the possibility that an object exists in this anchor, off-set value of BBOX refers to the off-set between the anchor BBOX and the ground truth BBOX, and the probability vector of all categories refers to the possibility that the object in this anchor belongs to each category. The number of anchors depends on the size of the feature map, in which each point has three anchors with different shapes predetermined by hyper-parameters. We have trained and tested the state-of-the-art object detectors on the TT100k dataset [30]. YOLO v5 achieves the best performance in terms of real time and accuracy. Therefore, we chose the YOLO v5 based TSR system as the attack target to evaluate our proposed physical adversarial attack.

### B. Physical Adversarial Examples

Most of the prior research either focused on the digital adversarial attacks, or launched physical adversarial attacks against the image classifiers [31]. However, successfully attacking a few static images cannot threaten the object detectors dealing with video streaming. There is no practical and satisfactory attempt to physically attack the object detectors due to three major challenges, *i.e.*, cross-domain conversion, image transformation, and limited capability.

*1) Cross-Domain Conversion:* As illustrated in Fig.2, AEs are generated in the digital domain, then printed into the physical domain, and finally re-taken by the camera back to the digital domain. The experimental results show that the digital-physical-digital conversion significantly degrades the adversarial toxicity of AEs [32]. To address this challenge, Jan *et al.* [33] presented an image-to-image translation network to simulate the digital-physical conversation. A conditional Generative Adversarial Network (cGAN) [34] is used to learn the digital-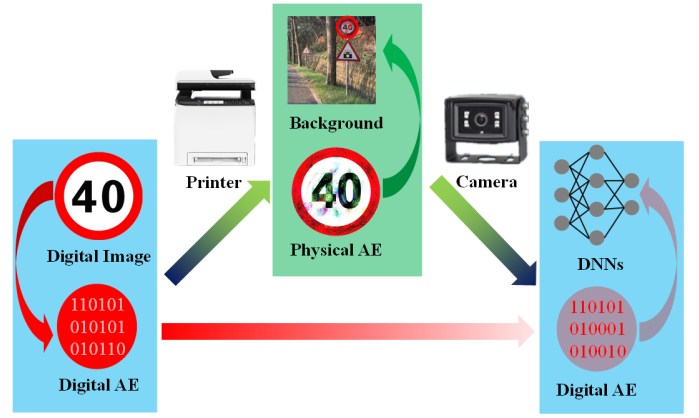physical conversion for generating a synthetic physical image, then it is served to produce adversarial noises in AEs generation. The cGANs model needs to be trained with a set of paired static images. Thus, the robustness of the physical AEs is improved when attacking several image classifiers.

*2) Image Transformation:* Different distances, angles, and illuminations will result in image transformations that impact the robustness of the physical AEs. AEs generated through the L-BFGS attack [5], fast gradient sign method [7], and the C&W attack [6] often lose their adversarial nature once subjected to minor transformations [35, 36]. To address this challenge, Athalye *et al.* [12] introduced the Expectation Over Transformation (EOT). EOT uses a chosen distribution of transformation functions and constrains the expected effective distance between AEs and the input images. EOT is able to generate 3D printed AEs which remain adversarial under a range of conditions. Within the framework of EOT, the physical AEs successfully deceive a standard image classifier.

*3) Limited Capability:* Further difficulty comes from the fact that the DNNs model is usually only a component in the whole computer vision system. For most applications, the attackers can neither get access to data inside the system nor obtain the parameters of the DNNs model. Instead, they can only manipulate objects in the physical environment. Therefore, the limitation of the attacker's capability limits the threat of the physical AEs. NaturalAE [37] proposed in 2021 used the natural AEs generated under white-box settings to attack the black-box models. However, the attack success rate is low, and the attack range is only $5m$. For a moving vehicle, the attack is only effective within less than half a second.

## III. THE PROPOSED PHYSICAL AE PIPELINE

In this section, we will elaborate on our proposed physical AEs pipeline. We first give the threat model with the goal and capabilities of the attacker. Then we present the 3-step AE attack pipeline which addresses the challenges in the previous section. Finally, we describe each of the three steps and show in detail the AEs generation algorithms for four attacks: Hiding Attack (HA), Appearance Attack (AA), Non-Target Attack (NTA), and Target Attack (TA).
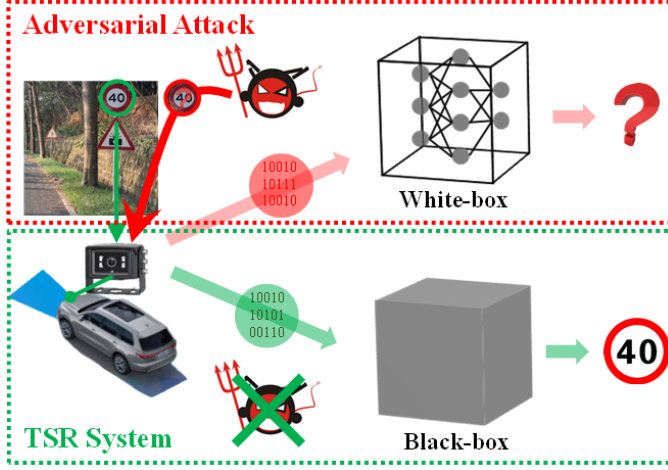
Fig. 3. Threat model of physical adversarial attack: the TSR system is a black box and the attacker does not have access to it.

## A. Threat Model

The attacker's goal is to mislead the victim vehicle's TSR system to incorrect classification of a given traffic sign. To this end, the attack can be further modeled as Hiding Attack (HA) where the attacker wants to hide the traffic sign, Appearance Attack (AA) where the attacker wants to create an object which can be recognized as a specific sign, Non-Target Attack (NTA) where the attacker simply wants the classification result to be incorrect, and Target Attack (TA) where the attacker wants the misclassified result to be a specific sign that is not the original traffic sign. Formally speaking, let $\boldsymbol{f}(\cdot)$ be the image classifier used in the TSR system, for a given input image $x$ whose category label is $y$, the attacker wants the TSR system to output a label that is different from $y$.

We adopt (i) the black-box TSR model where the attacker does not know the implementation details of the TSR system such as its structure and parameters. Furthermore, we assume that (ii) the attacker does not have physical access to the in-vehicle networks to analyze or modify the data in the TSR system. However, the attacker (iii) has access to the traffic sign which he can physically modify. Under these assumptions, the attacker's goal can be formulated as, generate an AE $x' = x + \delta$ for the input image $x$, where $\delta$ is a perturbation to $x$, such that:

$$\boldsymbol{f}(x') \neq y \qquad (1)$$

We believe that this is the most natural and strongest model for the attack scenario where an attacker wants to fool the TSR system of a victim vehicle. Any prior knowledge of the victim vehicle's information such as the model of the TSR, the object detector it uses, the camera(s) equipped in the vehicle, and the driving status (*e.g.* speed, weather, road condition) might provide additional help to the attacker and make the attack easier.

It is important to clarify that the previous AEs in the digital domain cannot be applied to this threat model. First, adversarial attack assumes a white-box model of the target DNN, the TSR system is a black-box and the parameters and structure of the target DNN models used in the TSR system

(assumption i) are not available to the attacker. So the attacker will not be able to generate digital AEs directly (see the top of Fig.3). Second, even if the attacker manages to create some digital AEs, because he does not have access to the in-vehicle networks (assumption ii), the attacker will not be able to inject such digital AEs directly into the TSR system (see bottom of Fig.3). However, the attacker can modify the traffic sign and hope the victim vehicle's TSR system will misclassify it.

## B. Attack Pipeline

As illustrated in Fig.4, there are three major steps in the proposed adversarial attack pipeline.

**Step 1.** As shown in Fig.4(a), in order to improve the robustness of the physical AEs, we extend the distribution of image transformations to simulate the changes of distance, angle, and illumination in the real world. The transformed images will be embedded into the background as the foreground to simulate the perspective of the object detectors.

**Step 2.** As shown in Fig.4(b), The BBOXes are associated with $x$, and can be extracted from the prediction results of YOLO v5. Single BBOX (S-BBOX) filter and Multiple BBOX (M-BBOX) filter are designed to obtain BBOXes for efficient perturbation training.

**Step 3.** As shown in Fig.4(c), each BBOX contains three types of information, *i.e.*, the probability of containing $x$ or $x'$, the probability vector for each category, and the BBOXes position offset. With the first two pieces of information, four loss functions are presented to generate AEs for aforementioned four attack vectors respectively.

## C. Image Preprocessing

The conventional AEs generation algorithms adopted the image transformations in EOT [12], including rotation, perspective, brightness, contrast, saturation, hue, and Gaussian noise. Furthermore, we extend distribution with the blur transformation to simulate conditions such as camera shake, and the resolution transformation that improves the robustness of AEs for varying distance to the camera. In step 1, we address the image transformation by setting:

$$x_t = \boldsymbol{t}(\boldsymbol{M}(x')) \qquad (2)$$

where $\boldsymbol{M}(\cdot)$ is the mask function to constrain the area where the perturbation is added. $\boldsymbol{t}(\cdot)$ is a transformation vector randomly selected from distribution of image transformations $T$.

We have attempted to fix the center point of the foreground at any position in the entire background. However, since part of the foreground is out of the background, it will make the perturbation training difficult to converge. Considering the fact that the camera will capture the complete traffic signs in most cases, we keep a certain distance between the center point and the edge of the background to guarantee that the entire foreground can be processed in the AEs generation algorithm. The image embedding algorithm is shown in Algorithm 1. During the optimization procedure, each transformed image should be embedded in different traffic backgrounds to imitate the perspective of the object detectors. $\boldsymbol{random}()$ can output
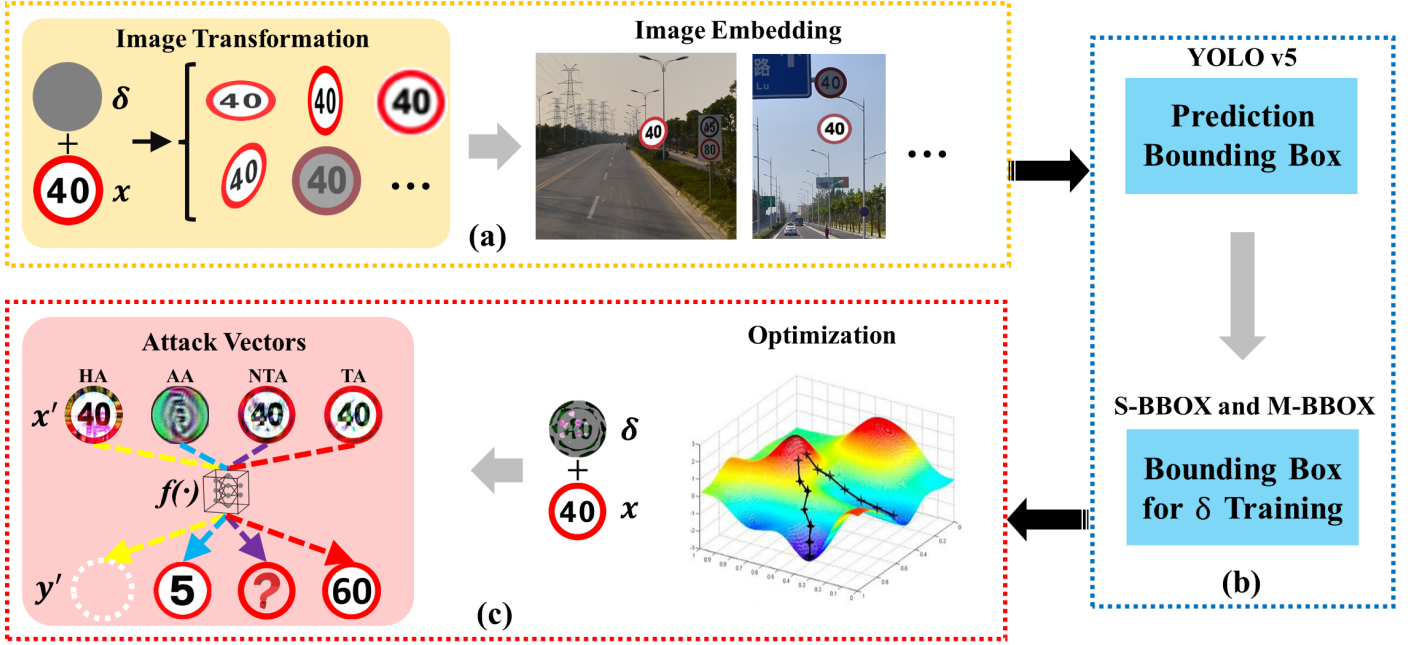
4

Fig. 4. Physical adversarial attack pipeline. (a) Image transformation and embedding. (b) BBOX filter. (c) Perturbation training for four attack vectors.

---

**Algorithm 1:** Image Embedding

**Input:** $r_1, r_2, r_3$: three ratios of foreground;
$\quad\quad\quad$ $i$: scale of background image.
**Output:** $bbox_{real}$: bounding box information of
$\quad\quad\quad\quad$ foreground.

1 **if** $\boldsymbol{random}() > 0.8$ **then**
2 $\quad$ $m = r_2 - r_1$
3 $\quad$ $a = r_1$
4 **else**
5 $\quad$ $m = r_3 - r_2$
6 $\quad$ $a = r_2$
7 $w_{real} = h_{real} = \boldsymbol{random}() \times m + a$
8 $s_f = \frac{3}{i}$
9 $p_x = \boldsymbol{random}() \times (1 - w_{real} - 2 \times s_f) + 0.5 \times w_{real} + s_f$
$\quad$ // avoid clinging the edge
10 $p_y = \boldsymbol{random}() \times (1 - h_{real} - 2 \times s_f) + 0.5 \times h_{real} + s_f$
11 $bbox_{real} = [p_x, p_y, h_{real}, w_{real}] \times i$

a random value in $[0,1]^{\mathbb{R}}$ every time it is called. $(p_x, p_y)$, $h_{real}$ and $w_{real}$ are the center point, height and width of the foreground respectively. $s_f$ is a position factor to prevent the foreground from clinging to the edge of background or even out of background. To simulate the real scale change of foreground in the background, we choose two different scale ranges corresponding to the big object and small object to generate the height and width of each BBOX. $r_1$, $r_2$, and $r_3$ indicate the ratio of the transformed images embedded in the background. The higher the ratio, the larger the object. We found that when the object is small, the direction of gradient update is hard to search and the model may not converge. However, if the object is far from the camera, lowering the ratio will enhance the effectiveness of AEs. Our empirical study shows that the best results come from when we choose

the ratio in [0.01, 0.1] with probability 20% and in [0.1, 0.5] with 80%. Therefore, $(r_1, r_2, r_3)$ is set to $(0.01, 0.1, 0.5)$ in our experiment.

### D. Bounding Box Filter

As mentioned before, YOLO v5 detects the object according to the feature maps of three different scales. For example, if the size of the input image is $640 \times 640$, the sizes of the three output feature maps are $20 \times 20$, $40 \times 40$, and $80 \times 80$ respectively. Each pixel in the feature maps is fixed with three anchors of different sizes to get the intermediate results, including the probability when there is an object in each anchor $(Q, 1)$, the confidence vector of the object category $(Q, N)$, and the off-set value vector between the real object and the fixed anchor $(Q, 4)$. Therefore, there are a total of $Q = 3 \times (20 \times 20 + 40 \times 40 + 80 \times 80)$ prediction BBOXes, and each has $S = 1 + 4 + N$ prediction values. In total, the final output of YOLO v5 is a matrix with dimension $Q \times S$. However, most of the prediction results of YOLO v5 are useless for perturbation training because some BBOXes detect either only some parts of $x$ or other objects. Therefore, we need to filter those useless BBOXes to improve the efficiency of the perturbation training.

We first extract the prediction BBOXes $O_{bbox}$ as defined below:

$$O_{bbox} = \boldsymbol{f}(\boldsymbol{emb}(b, x_t, bbox_{real})) \quad (3)$$

where $b$ refers to the background image, $\boldsymbol{f}(\cdot)$ is the object detector (YOLO v5). $\boldsymbol{emb}(b, x_t, bbox_{real})$ is an embedding function, in which the foreground object $x_t$ is embedded into the background image $b$ according to $bbox_{real}$ (Algorithm 1).

Two $\boldsymbol{BF}(\cdot)$ methods, S-BBOX and M-BBOX, have been proposed to filter the prediction BBOXes of YOLO v5 and

---

**Algorithm 2:** M-BBOX

---

**Input:** $bbox_{pred}$: YOLO v5' s result metric;
$bbox_{real}$: BBOXes obtained at algorithm 1;
$bbox_{anchor}$: BBOXes of fixed anchors.

**Output:** $bbox_{training}$: BBOXes for the perturbation training

1   $bbox_{offset}, bbox_{conf} = \boldsymbol{extract}(bbox_{pred})$
2   $bbox = \boldsymbol{plus}(bbox_{offset}, bbox_{anchor})$
3   $IOUs = \boldsymbol{iou}(bbox, bbox_{real})$
4   **if** *attack is HA* **then**
5      $index = \boldsymbol{where}(IOUs > 0.5)$
6      $middle = bbox_{pred}[index]$
7      $middle_{conf} = bbox_{conf}[index]$
8      $index_{conf} = \boldsymbol{top}(middle_{conf}, k)$
9      $bbox_{training} = middle[index_{conf}]$
10   **else if** *attack is AA or NTA* **then**
11      $index = \boldsymbol{top}(IOUs, k)$
12      $bbox_{training} = bbox_{pred}[index]$

---

extract $k$ BBOXes for the perturbation training. S-BBOX is used for TA, while M-BBOX has two modes, one is used for HA and the other for AA and NTA.

Then, we decompose the extracted $k$ BBOXes and obtain the information for the perturbation training:

$$V, P = \boldsymbol{split}(\boldsymbol{BF}(O_{bbox})) \qquad (4)$$

where $\boldsymbol{split}(\cdot)$ is the matrix split function. The first part is the probability $P \in \mathbb{R}^k$ of each target box containing the object $x$ or $x'$, the second part is the confidence vector $V \in \mathbb{R}^{k \times N}$ of the object category. Finally, $V$ and $P$ are used in the four loss functions corresponding to four attack vectors to train the perturbation.

*1) S-BBOX:* The Intersection Over Union (IOU) value between two BBOXes is calculated as follows:

$$\boldsymbol{iou}(bbox_A, bbox_B) = \frac{\boldsymbol{area}(bbox_A) \cap \boldsymbol{area}(bbox_B)}{\boldsymbol{area}(bbox_A) \cup \boldsymbol{area}(bbox_B)} \quad (5)$$

S-BBOX first filters out most of the prediction BBOXes whose detection confidence is lower than the threshold (equal to the NMS threshold of YOLO v5). Then, S-BBOX calculates the IOU value between each prediction BBOX and $bbox_{real}$ by Equ. (5). The prediction BBOX with the highest IOU value is for the perturbation training ($k = 1$ in S-BBOX).

*2) M-BBOX:* The M-BBOX defined in Algorithm 2 has two modes, a hiding mode for HA, and a non-hiding mode for AA and NTA. $\boldsymbol{extract}(\cdot)$ extracts the $bbox_{offset}$ and $bbox_{conf}$ from the $bbox_{pred}$. $\boldsymbol{where}(\cdot)$ obtain the indexes that satisfy the confidence in the brackets. $\boldsymbol{top}(A, k)$ search the indexes of prior $k$ in the ranking of vector $A$. $\boldsymbol{plus}(\cdot)$ achieve that adding the coordinate value of two input BBOXes.

Hiding Mode: When launching HA, we first calculate the IOU values between all the prediction BBOXes and the $bbox_{real}$. Then those BBOXes with IOU greater than the threshold are filtered out (Experimental results show that 0.5

is the best threshold). Finally, the filtered BBOXes are sorted according to their object confidence, and the prior $k$ BBOXes in the ranking of confidence are used to train perturbation.

Non-Hiding Mode: When launching AA and NTA, we first calculate the IOU values as in hiding mode. Then we directly obtain the BBOXes in the top $k$ IOU values to train perturbation.

*E. Four Attack Vectors*

The final goal of the four attack vectors is to deceive the target DNN models like YOLO v5 as imperceptible as possible. The optimization function is defined by:

$$arg \min_{x, \delta \in \mathbb{R}^{h \times w \times \varepsilon}} \mathbb{E}_{b \sim B, t \sim T} \boldsymbol{loss}_*(V, P, y, y') \qquad (6)$$

where $B$ refers to the background collection. We use $L_2$ distance to constrain the size of perturbation to improve its imperceptibility. The $L_2$ distance loss function is:

$$\mathcal{L}_{dis} = \|x - x'\|_2^2 \qquad (7)$$

*1) Hiding Attack:* The goal of HA is to make the object detector fail to find the object. The perturbation needs to eliminate the features of $x$ as a traffic sign by the loss function as below:

$$\boldsymbol{loss_H} = \frac{c}{nk} \sum_{i=1}^{n} \sum_{j=1}^{k} (\frac{1}{1 - P_i^j}) + \mathcal{L}_{dis} \qquad (8)$$

*2) Appearance Attack:* The main goal of AA is to make the object detector misrecognize AEs that cannot be recognized by human eyes as desired objects. This attack is to train perturbation on a blank image that can be recognized by the object detector as a specified category. Thus, the minimum Euclidean distance between AEs and the blank image is not required in AA. The loss function for AA is given as:

$$\boldsymbol{loss_A} = \frac{c}{nk} \sum_{i=1}^{n} \sum_{j=1}^{k} (\frac{1}{P_i^j \times V_i^{j,y'}}) \qquad (9)$$

*3) Non-Target Attack:* The goal of NTA is to make the object detector misrecognize AEs that belong to a certain category as some other categories. The perturbation needs to eliminate $x$'s features of the correct category, but retain $x$'s features of a traffic sign. The loss function for NTA can be defined as:

$$\boldsymbol{loss_{NT}} = \frac{c}{nk} \sum_{i=1}^{n} \sum_{j=1}^{k} (\frac{1}{P_i^j} + \frac{1}{1 - V_i^{j,y}}) + \mathcal{L}_{dis} \qquad (10)$$

Fig. 5. Impact of $c$ on imperceptibility of AEs. (a) $c$ is set to a small value. (b) $c$ is set to a medium value. (c) $c$ is set to a large value.

*4) Target Attack:* The goal of TA is to make the object detector misrecognize AEs that belong to a certain category as the target category $y'$. The perturbation needs to not only change $x$'s features as NTA, but also add the features of $y'$. Thus, it will be tricky to design the loss function in TA. The loss function can be defined as:

$$loss_T = \frac{c}{nk} \sum_{i=1}^{n} \sum_{j=1}^{k} \left( \frac{1}{P_i^j} + \frac{1}{V_i^{j,y'}} + \sum_{z=1,z\neq y}^{N} \frac{1}{1 - V_i^{j,z}} \right) + \mathcal{L}_{dis}$$
(11)

where $N$ is the number of categories that YOLO v5 can detect. The adjustable parameter $k$ is the number of target boxes that YOLO v5 extracts from each sample. In each loss function, $c$ needs to be adjusted in each training phase according to the actual environmental conditions. $n$ is the number of samples used for training. As shown in Fig.5, the larger the value of $c$ is, the greater the perturbation will be added, and hence the more robust the generated AEs will be.

## IV. EVALUATION

We launch the four attack vectors (*i.e.*, HA, AA, NTA, and TA) against the state-of-the-art object detectors and a brand-new vehicle. Considering that the TSR system in the vehicle only recognizes the speed limit signs, we focus on physically attacking the speed limit of $40km/h$ (limit 40 for short), which is generally located on the exit ramp of the freeway thus is more life-threatening. Due to the space limit, we cannot present all the experimental results in this paper. More than 1,000 video clips and more evaluation results are uploaded on our demo website: https://seczone.cn/contents/422/1024.html.

### A. Experimental Setup

AEs are generated on the server equipped with two Intel Xeon-E5-2680 V4 CPUs, four NVIDIA GTX3080 GPUs, and

64GB RECC DDR4-2400MHz memory. Then we fabricate the physical AEs in accordance with the production specifications of the traffic signs, including the material and size, *etc*. We conduct multiple sets of outdoor experiments under a variety of environmental conditions. The distances range from $0m$ to $30m$, the angles range from $-60°$ to $60°$, and the illuminations range from sunny days, cloudy days to nights. The videos are captured by the built-in cameras of the Huawei nova7 with 64 Mega pixels and aperture of f/1.8, and the Samsung S10 with 16 Mega pixels and aperture of f/2.4. The effectiveness and robustness of the physical AEs are evaluated by analyzing each image frame in the video streaming. In order to present the threat of the physical adversarial attack more comprehensively, we not only count the attack success rate but also measure the impact of AEs on the detection confidence.

*1) Dataset:* The YOLO v5 based TSR system is trained with the TT100k dataset [30], which is composed of $2048 \times 2048$ images. The TT100k dataset declares that there are 221 categories containing 100,000 images of 30,000 traffic sign samples. However, we only manage to obtain 16,823 images with 26,337 traffic sign samples, 6,107 images for training, 3,073 for testing, and 7,643 other images from its website. We use the original TT100K dataset to train the object detectors first. The detectors will not converge when the high-resolution images are directly resized into the $640 \times 640$ input images for training. Thus, we divide each $2048 \times 2048$ image into 16 $640 \times 640$ images to train the TSR system. Besides, we find that only 150 categories have data in the TT100k dataset, and among the nonempty categories, half of them have less than 10 images. Such unbalanced distribution of training data will severely impact the performance of the DNNs models [38]. After many training experiments, we screened out the traffic sign categories with top 50 data volume to form the new dataset for our experiments. The new TT100k dataset has 50 categories of 21,881 images with $640 \times 640$, 14474 images for training, 7,407 images for testing, and contains 40,550 traffic sign samples.

Table I shows the performance of the object detectors on the new dataset. Among the six detectors, Faster R-CNN [22] is a two-stage detection, others are all one-stage target detectors; CenterNet [41] is an anchor-free detector, others are all anchor-base detectors. The one-stage detectors, YOLO v3 [25], YOLO v5 [26], and SSD [39], have high detection accuracy, even outperforming the two-stage detector Faster R-CNN [22]. CenterNet [41] has lower detection accuracy than most anchor-base detectors on the new dataset. In general, YOLO v5 [26] performs the best on the new dataset.

*2) Hyper-Parameters Setup:* Stochastic Gradient Descent (SGD) algorithm [42] is utilized for the optimization of the
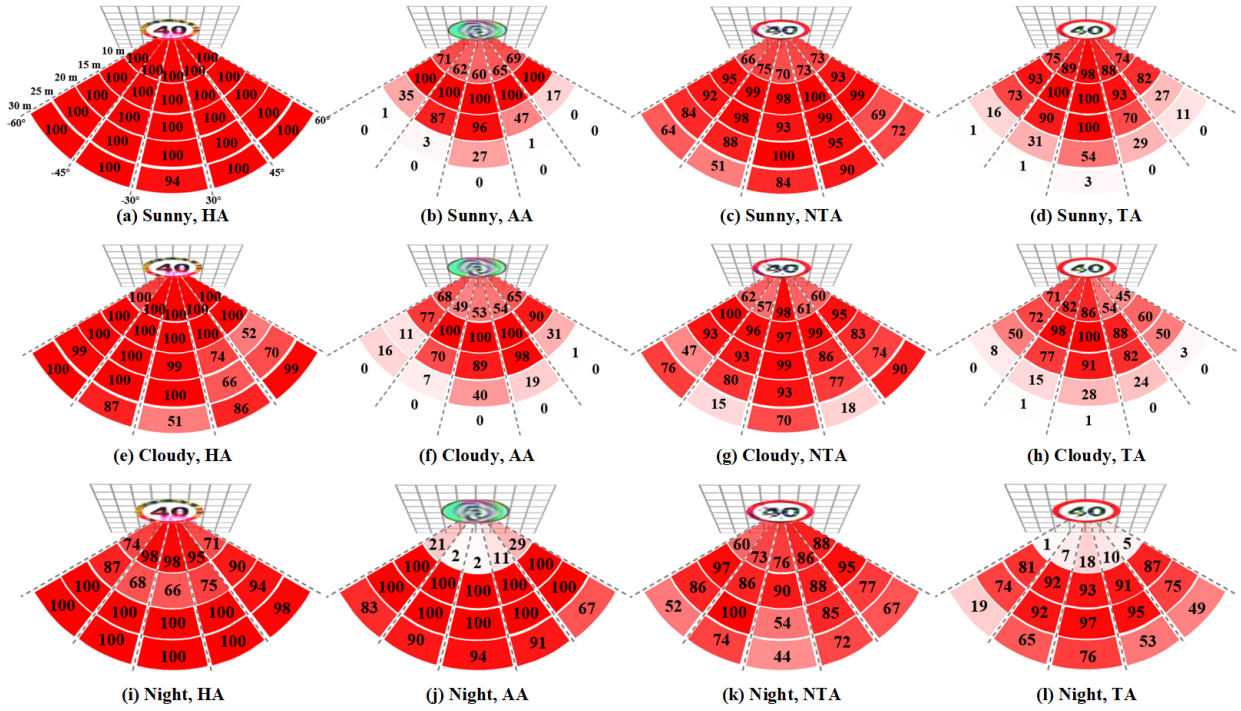
Fig. 6. Attack success rate of four attack vectors under a variety of environmental conditions. (a) HA on a sunny day. (b) AA on a sunny day. (c) NTA on a sunny day. (d) TA on a sunny day. (e) HA on a cloudy day. (f) AA on a cloudy day. (g) NTA on a cloudy day. (h) TA on a cloudy day. (i) HA at night. (j) AA at night. (k) NTA at night. (l) TA at night.

loss functions, in which the transformation batch is set to 16, and the size of the input image and the background image is set to $640 \times 640$. For the four attack vectors, we set different experimental hyperparameters as follows:

**Hiding Attack.** HA aims to hide limit 40 in the background. In HA, the $c$ is set to $1e + 2$, the learning rate is set to 0.1, and k is set to 10.

**Appearance Attack.** AA aims to create a traffic sign which can be recognized as limit 5 by YOLO v5 from a blank image. In AA, $x$ is a randomly sampled metric with $640 \times 640$. $c$ is set to $1e + 5$, the learning rate is set to 0.2, and k is set to 1.

**Non-target Attack.** NTA aims to make YOLO v5 recognize the limit 40 as other categories. In NTA, $c$ is set to $1e+5$, the learning rate is set to 0.1, and k is set to 3.

**Target Attack.** TA aims to make YOLO v5 recognize the limit 40 as the limit 60. In TA, the $c$ is set to $1e + 3$, the learning rate is set to 0.1, and k is set to 1.

*3) Evaluation Metrics:*

$$N_s = \begin{cases} \mathbf{Z}(x) > th \cap \mathbf{Z}(x') < th, & HA \\ \mathbf{g}(V_{x'}) = y', & AA \\ \mathbf{g}(V_x) = y \cap \mathbf{g}(V_{x'}) \neq y, & NTA \\ \mathbf{g}(V_x) = y \cap \mathbf{g}(V_{x'}) = y', & TA \end{cases}$$

$$N_a = \begin{cases} 1, & HA\&AA \\ \mathbf{g}(V_x) = y, & TA\&NTA \end{cases}$$

(12)

The attack success rate is defined as $R_s = \sum N_s / \sum N_a \times 100\%$. As shown in Equ. (12), $\mathbf{g}(\cdot)$ denotes the $\mathbf{argmax}(\cdot)$ function that outputs the index of the max value in the vector.
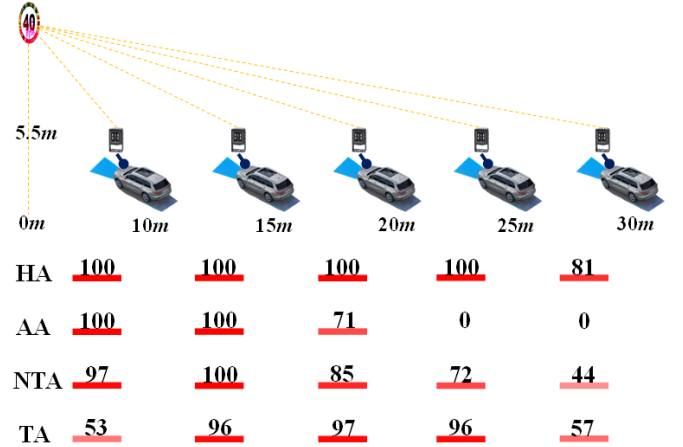


Fig. 7. Impact of height on attack effectiveness.

In HA and AA, $N_a$ is every image frame captured by the camera, while in NTA and TA, $N_a$ is the image frame in which the object is correctly detected. $N_s$ depends on the attack goals of the four attack vectors. In HA, a successful attack means that the probability of detecting $x$ ($\mathbf{Z}(x) = \mathbf{max}(V_x) \times P_x$) is greater than the threshold $th = 0.25$, and the probability of detecting $x'$ ($\mathbf{Z}(x') = \mathbf{max}(V_{x'}) \times P_{x'}$) is less than the threshold. In AA, a successful attack means $x'$ is detected as the target category $y'$. In NTA, a successful attack means $x$ is detected as $y$, but $x'$ is not detected as the original category $y$. In TA, a successful attack means $x$ is detected as $y$, and $x'$ is detected as the target category $y'$.
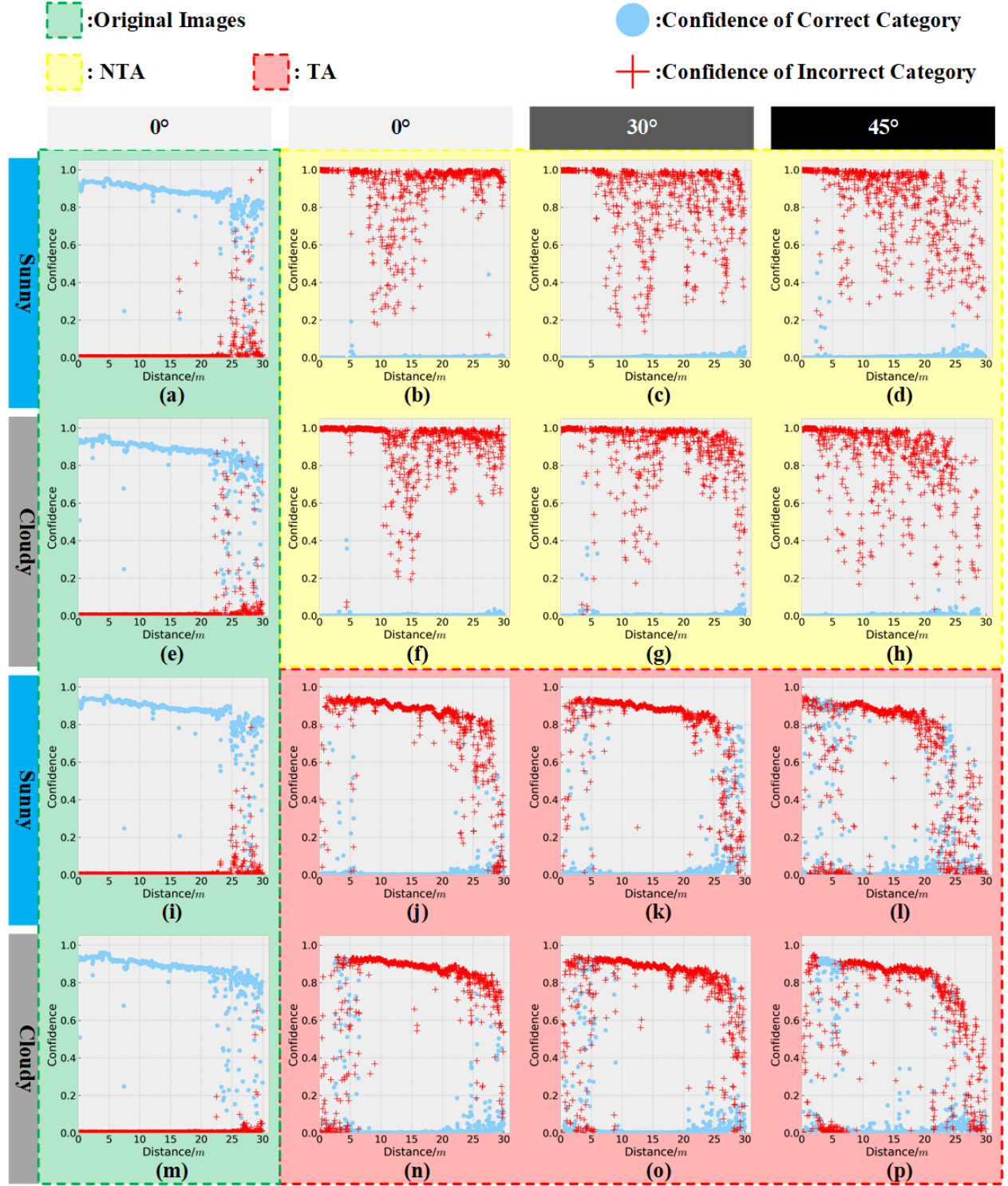
Fig. 8. Confidences of NTA and TA under a variety of environmental conditions. (a) NTA with original image on a sunny day. (b) NTA at $0°$ on a sunny day. (c) NTA at $30°$ on a sunny day. (d) NTA at $45°$ on a sunny day. (e) NTA with original image on a cloudy day. (f) NTA at $0°$ on a cloudy day. (g) NTA at $30°$ on a cloudy day. (h) NTA at $45°$ on a cloudy day. (i) TA with original image on a sunny day. (j) TA at $0°$ on a sunny day. (k) TA at $30°$ on a sunny day. (l) TA at $45°$ on a sunny day. (m) TA with original image on a cloudy day. (n) TA at $0°$ on a cloudy day. (o) TA at $30°$ on a cloudy day. (p) TA at $45°$ on a cloudy day.

## B. Evaluation of Four Attack Vectors

*1) Attack Success Rate:* As shown in Fig.6, the four attack vectors are evaluated in terms of the attack success rate $R_s$ under a variety of environmental conditions. In this set of experiments, we fixed the height of the traffic signs to about $1.5m$. The distances are divided into five regions, *i.e.*, $[0m, 10m]$, $[10m, 15m]$, $[15m, 20m]$, $[20m, 25m]$, $[25m, 30m]$. The video clips are recorded from far to close with a constant speed in each region at the angle of $[-60°, -45°]$, $[-45°, -30°]$, $[-30°, 30°]$, $[30°, 45°]$, and $[45°, 60°]$, respectively. For each $R_s$ in Fig.6, more than 200 image frames are captured and processed (Fig.13 and Fig.14 in appendix). Since the visibility is low at night, we only conduct experiments within $25m$. The depth of red represents the value of $R_s$, the darker the color indicates the higher the attack success rate.

For HA, we achieve the best attack effectiveness on a sunny day. On a cloudy day, in the region of $[15m, 20m]$ at $[45°, 60°]$ and in the region of $[25m, 30m]$ at $[-30°, 30°]$, $R_s$ degrades to 50%. At night, in the region of $[10m, 15m]$ at $[-45°, -30°]$ and $[-30°, 30°]$, $R_s$ degrades below 70%. However, HA has a high $R_s$ at $10m$ away in the dark environment. The experimental results indicated that reflection at night has a greater impact on HA.

For AA, we achieve the best attack effectiveness at night, especially in the region of $[10m, 20m]$. However, in the dark environment within $10m$, $R_s$ is very low. On sunny and cloudy days, $R_s$ is high in the region of $[10m, 15m]$, but AA almost failed at $25m$ away. Thus, $R_s$ is relatively higher in the darker light, while both larger angles and distances will significantly reduce $R_s$. The experimental results indicated the huge impact of reflection at night on AA too.

For NTA, since the strong reflection caused by direct lighting covers the perturbation features, $R_s$ degrades below 60% in the region of $[15m, 20m]$ at $[-30°, 30°]$ at night. On a cloudy day, $R_s$ degrades below 20% in the region of $[25m, 30m]$ at $[-45°, -30°]$ and $[30°, 45°]$. The results indicate that both darker light and larger distances reduce $R_s$, while the impact of different angles on $R_s$ does not show a strong regularity.

For TA, we achieve the best attack effectiveness in the region of $[10m, 20m]$ at $[-45°, 45°]$. At night, TA almost fails within $10m$ in the dark environment. On a sunny day and a cloudy day, $R_s$ degrades severely as the distance increases at $20m$ away. The results indicate that large angles reduce $R_s$ in TA when the distance is either very close or very far; while the impact of different illuminations on $R_s$ does not show a strong regularity.

Then we increased the height of the traffic sign to about $5.5m$ and conducted the four attack vectors on a sunny day at $0°$ as shown in Fig.11 in Appendix.

The experimental results in Fig.7 indicate that the height of AEs has little impact on HA within $25m$, but $R_s$ drops by 13% in the region of $[25m, 30m]$. For AA, $R_s$ rises from 60% to 100% in the region of $[0m, 10m]$ when the height of AEs rises from $1.5m$ to $5.5m$, but $R_s$ drops by 25% in the region of $[15m, 20m]$. Thus, $R_s$ in AA increases with height in the closer region, while at $20m$ away, $R_s$ in AA decreases as the height of AEs increases. Similarly, NTA has higher $R_s$ in the

| Attack Vector | Number of image frames | $R_s$/% |
|---|---|---|
| HA | 824 | 96.48 |
| AA | 630 | 60.48 |
| NTA | 525 | 90.48 |
| TA | 645 | 92.87 |

| | Faster R-CNN[22] | SSD[39] | RetinaNet[40] | CenterNet[41] |
|---|---|---|---|---|
| HA | 95.02 | 95.40 | 71.02 | 97.64 |
| AA | 54.73 | 28.23 | 11.94 | 19.65 |
| NTA | 99.78 | 100 | 52.18 | 47.8 |
| TA | 58.62 | 46.38 | 2.03 | 0 |
| mAP-s | 61.5 | 64.6 | 48.5 | 55.3 |
| mAP-m | 77.7 | 79.3 | 77.2 | 77.9 |
| mAP-l | 82.5 | 83.7 | 84 | 85.7 |

region of $[0m, 20m]$ and lower $R_s$ at $20m$ away when the height of AEs rises from $1.5m$ to $5.5m$. For TA, $R_s$ drops by 45% in the region of $[0m, 10m]$, but rises from 54% to 96% in the region of $[20m, 25m]$ and rises from 3% to 57% in the region of $[25m, 30m]$ when the height of AEs rises from $1.5m$ to $5.5m$.

We also test the four attack vectors in a vehicle with a speed of $[20km/h, 30km/h]$ as shown in Fig.12 in Appendix. The test is conducted in the region of $[0m, 30m]$ at $[-30°, 30°]$ on a sunny day. The attack effectiveness is listed in Table II. The average $R_s$s in HA, NTA, and TA are all over 90%. The results of the real-road driving test for AA are consistent with the results in Fig.6. The experimental results in Table II show that our generated AEs are robust in the real-road driving test.

*2) Confidence:* In order to further study the effectiveness and robustness of NTA and TA, we measure the detection confidence of each image frame. We fix the height of the traffic signs to about $1.5m$ and choose the video clips recorded on a sunny day and a cloudy day in the region of $[0m, 30m]$ at $0°$, $30°$, and $45°$, respectively. As shown in Fig.8, the blue circle presents confidence of the correct category $y$. It needs to be emphasized that in NTA, the red cross represents the highest confidence of all the incorrect categories, while in TA, it represents the confidence of the target category $y'$.

Fig.8(a,e,i,m) in the leftmost column show confidence in each original input image. Limited by the resolution of the camera, the detection performance of the original images drops significantly at $25m$ away. In several image frames, the confidence of the correct category $y$ drops severely, while the confidence of the incorrect categories increases.

By comparing the results of NTA on a sunny day in Fig.8(b,c,d) and the results of NTA on a cloudy day in Fig.8(f,g,h), we find that the illumination has little impact on the attack effectiveness of NTA. By comparing the results of TA on a sunny day in Fig.8(j,k,l) and the results of TA on a cloudy day in Fig.8(n,o,p), it indicates that the attack effectiveness of TA is better in a sunny day within $10m$. By
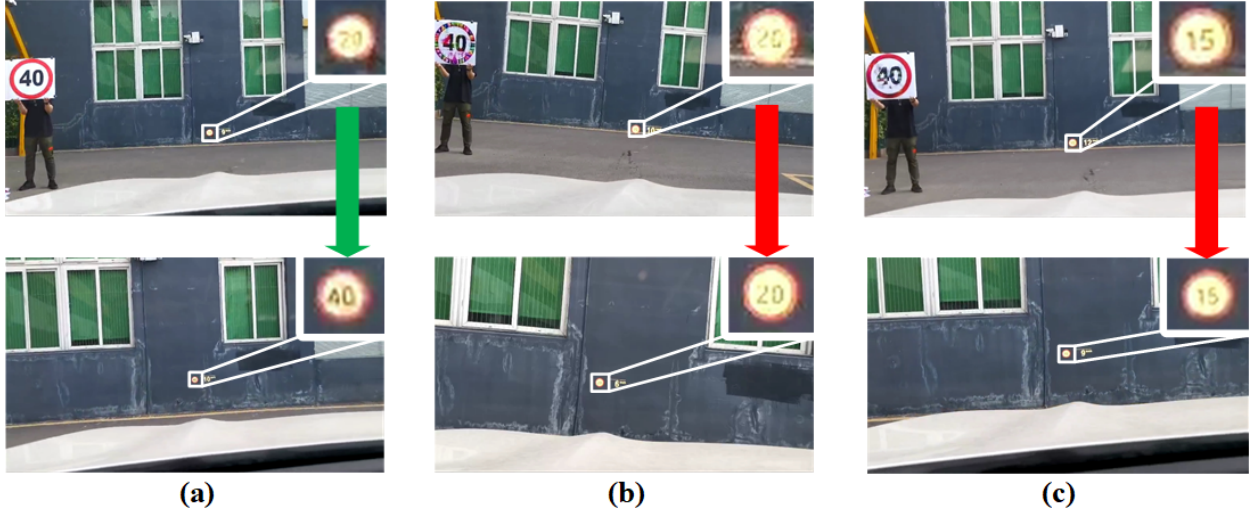
Fig. 9. HA and NTA against a brand-new vehicle. (a) The original speed limit 40. (b) HA. (c) NTA.

comparing the experimental results at $0°$ in Fig.8(b,f,j,n), $30°$ in Fig.8(c,g,k,o), and $45°$ in Fig.8(d,h,l,p), it indicates three phenomena. First, the large angle has a greater impact on TA than NTA. Second, at very close and far distances, the large angle significantly degrades the attack effectiveness of TA. Third, attack effectiveness is better in a sunny day than in a cloudy day.

*C. Attack Transferability*

*1) State-of-the-art Object Detectors:* We use AEs generated with YOLO v5 to attack the representative object detectors, including Faster R-CNN (two-state) [22], SSD [39] and RetinaNet [40] (one-stage and anchor-base), and CenterNet (one-stage and anchor-free) [41]. Those object detectors are trained on the same TT100k [30] dataset for fairness. The set of experiments is conducted in the region of $[0m, 30m]$ at $0°$. According to the experimental results in Table III, HA and NTA achieve higher $R_s$, especially on SSD and Faster R-CNN, while AA and TA on RetinaNet and CenterNet almost fail. We measure the performance of the four detectors when detecting large, medium and small objects, and find that RetinaNet and Centernet have low accuracy when detecting small objects. Therefore, low accuracy at long distances determines that they have low transferability in the region of $[20m, 30m]$, but they still have high accuracy at close distances, and thus have good transferability in the region of $[0m, 15m]$.

*2) Brand-new Vehicle:* The major reason why we generated adversarial limit 40 is that the TSR system of the brand-new vehicle can only recognize the speed limit signs. As shown in Fig.9(a), only after the vehicle passes the speed limit sign, the recognition result can be displayed on the screen. In HA and NTA, Fig.9(b,c) show that after the vehicle passes the adversarial limit 40, it does not display the correct limit 40, which means HA and NTA successfully fool the eyes of the vehicle.

*D. Comparison and Discussion*

There are few works focusing on the physical AEs against the object detectors. As listed in Table IV, We compare our proposed method with Zhao's method [18], ShapeShifter [31], and NaturalAE [37] from different dimensions. In general, our method has the following advantages. First of all, to the best of our knowledge, it is the first attempt to physically attack a production-grade object detector in a real-world vehicle and achieve significant effectiveness. Second, we take into account the impact of light reflection in the dark environment on robustness and conduct a complete set of experiments at night. Surprisingly, all the four attack vectors achieve almost 100% $R_s$ against the YOLO v5 based object detector at a certain distance and angle. Third, NTA and TA are implemented in the physical domain and exhibit satisfactory transferability against the state-of-the-art object detectors, which lays a foundation for studying the adversarial attacks against black-box models in both the digital and the physical domains.

Combining all the experimental results, we can draw a conclusion that the physical AEs for HA and NTA have better robustness and transferability against multiple object detectors, while AA and TA have high $R_s$ only on specific attack targets. The lower $R_s$ and poorer transferability are due to two-fold reasons. First, HA only needs to eliminate the features of all the categories in the perturbation training process, and NTA only needs to eliminate the features of the correct category $y$. But for AA and TA, the features of the target category $y'$ should be imitated, which is much more difficult than the feature elimination. In addition, TA also needs to eliminate the features of $y$ before the feature imitation. Therefore, the digital-physical-digital conversion has a greater impact on the AEs generated for AA and TA, which results in lower transferability. Second, the similarity of the two target models also determines the transferability of the physical AEs. The architectures of RetinaNet and CenterNet are quite different from YOLO v5. As for the TSR system of the brand-new vehicle, even the architecture of the DNNs model is unknown. The black-box setting is always the biggest obstacle for the adversarial attack in both digital and physical domains. In future work, we will launch the physical adversarial attack against the black-box DNNs model with the help of the side-channel attack [43].

| | Our Method | Zhao's Method [18] | ShapeShifter [31] | NaturalAE [37] |
|---|---|---|---|---|
| **Target Model** | YOLO v5 | YOLO v3 | Faster R-CNN | YOLO v2 |
| **Attack Vector** | HA, AA, NTA, TA | HA, AA | TA | TA |
| **Distance** | [0 m, 30 m] | [0 m, 25 m] | [0 m, 12 m] | [0 m, 5 m] |
| **Angle** | [-60,60] | [-60,60] | [0,60] | [-60,60] |
| **Illumination** | Sunny, Cloudy, Night | Indoor, Outdoor | Indoor, Outdoor | Indoor, Outdoor |
| **Height** | 1.5 m, 5.5 m | 1.5m | 1.5m | 1.5m |
| **Transferability** | Faster R-CNN, SSD, RetinaNet, CenterNet Brand-new vehicle (HA and NTA) | Faster R-CNN, SSD RFCN, Mask R-CNN | None | Faster R-CNN, SSD |



Fig. 10.    Adversarial Attacks.

## V.    RELATED WORK

We survey the landmark works related to the adversarial attacks in Fig.10, and discuss the AEs defense mechanisms to protect the DNNs models from the adversarial attacks.

### A. Adversarial Attacks

*1) Digital AEs:* In 2013, Szegedy *et al.* [5] discovered the existence of AEs in the digital domain and generated AEs using the box-constrained L-BFGS. Since then, there have been a lot of initial works [6–10] focusing on generating digital AEs. Goodfellow *et al.* [7] proposed the fast gradient sign method to generate AEs on the MNIST dataset. Carlini and Wagner [6] introduced three C&W attacks for the $L_0$, $L_2$, and $L_\infty$ distance metrics. The $L_0$ attack was the first published attack that caused the targeted misclassification on the ImageNet dataset. Although most of the digital AEs lose their adversarial nature in the physical environment [35, 36], these three classic methods are still used for generating physical AEs.

*2) Physical AEs:* Athalye *et al.* [12] proposed the Expectation Over Transformations (EOT) method that laid the cornerstone of the physical AEs generation algorithms. Evtimov *et al.* [44, 45] used EOT method to generate robust physical AEs to attack YOLO v2 based TSR systems. Sitawarin *et al.* [14] presented two novel attack vectors against the traffic sign classifier. Based on [14], Morgulis *et al.* [46] claimed to attack the real production-grade image classifiers for the first time. Chen *et al.* [31] physically attacked the fast R-CNN based TSR systems using large perturbations. Lovisotto *et al.* [47] used the light of a projector to generate short-lived adversarial perturbations in the indoor experiments. However, the physical AEs generated by these methods cannot remain robust under a variety of environmental conditions.

Until very recently, there were several attempts to improve the robustness of the physical AEs. Jan *et al.* [33] used an image-to-image translation network to simulate the digital-physical conversion for generating the physical AEs. This improvement performed well on several image classifiers with static input images, but it was not verified on the object detectors with video streaming. Xue *et al.* [37] proposed the natural and robust physical AEs against the object detectors. Since the generated AEs are effective within 5 meters, they pose almost no threat to a moving vehicle. Zhao *et al.* [18] presented the hiding attack and appearance attack on the YOLO v3 based object detector and achieved better robustness. Limited by the capability of the AEs generation algorithm, the imperceptibility of AEs is not properly addressed which makes the two attack vectors lack real-world threat to the production-grade TSR systems.

### B. Defense Mechanisms

Research has proposed several defense mechanisms to protect the DNNs models against AEs, which can be divided into three categories, *i.e.*, input image preprocessing, AEs detection, and model enhancement.

*1) Input Image Preprocessing:* Imperceptible requirements make AEs not robust enough to the external noises or data distortions [48]. Inspired by this opportunity, preprocessing the input image to eliminate the adversarial features could be a type of potential defense mechanism. In 2017, Das *et al.* [49] explored and demonstrated that the system's JPEG compression could be used as an effective preprocessing step in the classification pipeline to significantly reduce the impact of AEs. JPEG compression has ability to remove the high-frequency signal components inside the image, which helps eliminate malicious disturbances. Guo *et al.* [50] selected a small group of pixels and reconstructs the simplest image consistent with the selected pixels so that malicious disturbance in the image was removed. Similar to [50], Prakash *et al.* [51]
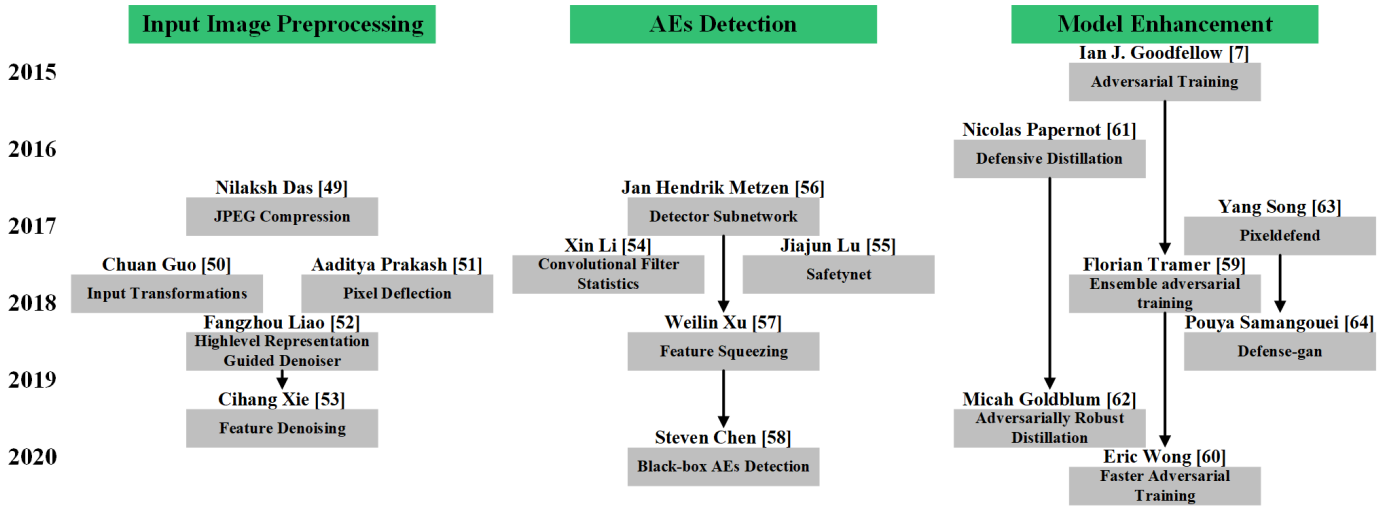
Fig. 11. Defense mechanisms against AEs. (a) Input image preprocessing. (b) AEs detection. (c) Model enhancement.

proposed a pixel deflection method that used semantic maps and randomization to select a small number of pixels, and then replace them with randomly selected neighboring pixels. This process would generate noise and needed a wavelet noise reduction filter. In order to solve this problem, Liao *et al.* [52] used U-net structure to modify the denoising autoencoder and proposed denoising U-net. However, the residual noise impact might still increase as the number of network layers increases. In 2019, Xie *et al.* [53] developed a new network architecture to improve robustness with a feature denoiser that combined confrontation training and graphics processing.

*2) AEs Detection:* It is also possible to directly detect the adversarial features and prevent AEs from entering DNNs models. Li *et al.* [54] used a cascade classifier to detect AEs effectively. Lu *et al.* [55] proposed SafetyNet to detect AEs based on that benign images and AEs have different ReLU activation patterns in the model. Metzen *et al.* [56] chose a small detector sub-network to enhance the robustness of the DNNs model. The sub-network would be trained on a binary classification task to distinguish AEs from benign images. On this basis, Xu *et al.* [57] proposed a feature compression method to detect AEs by comparing the prediction consistency of the original image and the compressed image. If the classification result of the original input image and the compressed input image were different, the input might be an AE. In 2020, Chen *et al.* [58] presented a defense mechanism to detect the process of AEs generation. By keeping the historical record of past queries, it could determine when a series of queries appeared to generate AEs.

*3) Model Enhancement:* Enhancing the DNNs models in the training phase is also a potential defense mechanism against adversarial attacks. In the next year after the AEs concept appeared, Goodfellow *et al.* [7] proposed to add AEs to the training set to make the trained DNNs model robust against specific types of AEs. Tramer *et al.* [59] believed that [7] was vulnerable to single-step attacks. They proposed an ensemble adversarial training method based on gradient masks, using AEs generated by other pre-trained classifiers to expand the training set. This ensemble adversarial training algorithm was more effective because it separated the training of the

model and the process of generating AEs. In 2020, Wong *et al.* [60] paid more attention to the cost of adversarial training, and proposed a fast adversarial training method. Defensive distillation was proposed by Papernot *et al.* [61] in 2016. The basic idea was to transfer knowledge of complex networks to simple networks by modifying the network structure and optimization items to prevent the model from fitting too closely to normal samples. However, the network needed to be retrained and was usually only effective against AEs considered in the training process. In 2020, Goldblum *et al.* [62] studied how the adversarial robustness is transferred from the teacher DNNs model to the student DNNs model in the knowledge distillation process, and introduced adversarial robust distillation to refine the robustness to the student DNNs. Another robustness enhancement method was to use a generative model to project potential AEs onto a benign dataset, and then classify them. Among them, the PixelDefend method proposed by Song *et al.* [63] used the PixelCNN generative model, and the Defence-GAN method proposed by Samangouei *et al.* [64] used a generative adversarial network structure.

## VI. Conclusion

In this paper, we present a systematic pipeline to generate the physical AEs against the object detectors. In order to improve the robustness of AEs in the physical domain, we extend the distribution of image transformations, design the S-BBOX filter and the M-BBOX filter, and modify the four loss functions for the four attack vectors respectively. We conduct extensive experiments under a variety of environmental conditions, *i.e.*, the distance varies from $0m$ to $30m$, the angle varies from $-60°$ to $60°$, the illumination varies from sunny day to cloudy day to night. The experimental results show that HA, NTA, and TA achieve a success rate of more than 90% in real-world driving tests against the YOLO v5 based TSR system. The generated AEs exhibit high transferability against the other state-of-the-art object detectors. Furthermore, HA and NTA successfully fooled the TSR system of a brand-new vehicle, which is a life-threatening case for autonomous vehicles.

REFERENCES

[1] R. Singh, A. Agarwal, M. Singh, S. Nagpal, and M. Vatsa, "On the robustness of face recognition algorithms against attacks and bias," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 09, 2020, pp. 13 583–13 589.

[2] R. J. S. Raj, S. J. Shobana, I. V. Pustokhina, D. A. Pustokhin, D. Gupta, and K. Shankar, "Optimal feature selection-based medical image classification using deep learning model in internet of medical things," *IEEE Access*, vol. 8, pp. 58 006–58 017, 2020.

[3] S. Lan, Z. Ren, Y. Wu, L. S. Davis, and G. Hua, "Saccadenet: A fast and accurate object detector," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 397–10 406.

[4] J. Zhang and C. Li, "Adversarial examples: Opportunities and challenges," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 7, pp. 2578–2593, 2019.

[5] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013.

[6] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 ieee symposium on security and privacy (sp)*. IEEE, 2017, pp. 39–57.

[7] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[8] N. Papernot, P. McDaniel, S. Jha, M. Fredrikson, Z. B. Celik, and A. Swami, "The limitations of deep learning in adversarial settings," in *2016 IEEE European symposium on security and privacy (EuroS&P)*. IEEE, 2016, pp. 372–387.

[9] J. Kos, I. Fischer, and D. Song, "Adversarial examples for generative models," in *2018 ieee security and privacy workshops (spw)*. IEEE, 2018, pp. 36–42.

[10] X. Yuan, Y. Chen, Y. Zhao, Y. Long, X. Liu, K. Chen, S. Zhang, H. Huang, X. Wang, and C. A. Gunter, "Commandersong: A systematic approach for practical adversarial voice recognition," in *27th {USENIX} Security Symposium ({USENIX} Security 18)*, 2018, pp. 49–64.

[11] A. Kurakin, I. Goodfellow, S. Bengio *et al.*, "Adversarial examples in the physical world," 2016.

[12] A. Athalye, L. Engstrom, A. Ilyas, and K. Kwok, "Synthesizing robust adversarial examples," in *International conference on machine learning*. PMLR, 2018, pp. 284–293.

[13] I. Evtimov, K. Eykholt, E. Fernandes, T. Kohno, B. Li, A. Prakash, A. Rahmati, and D. Song, "Robust physical-world attacks on machine learning models," *arXiv preprint arXiv:1707.08945*, vol. 2, no. 3, p. 4, 2017.

[14] C. Sitawarin, A. N. Bhagoji, A. Mosenia, M. Chiang, and P. Mittal, "Darts: Deceiving autonomous cars with toxic signs," *arXiv preprint arXiv:1802.06430*, 2018.

[15] Y. Jin, Y. Fu, W. Wang, J. Guo, C. Ren, and X. Xiang, "Multi-feature fusion and enhancement single shot detector for traffic sign recognition," *IEEE Access*, vol. 8, pp. 38 931–38 940, 2020.

[16] Z. Wang, Y. Wu, L. Yang, A. Thirunavukarasu, C. Evison, and Y. Zhao, "Fast personal protective equipment detection for real construction sites using deep learning approaches," *Sensors*, vol. 21, no. 10, p. 3478, 2021.

[17] L. Schmarje, M. Santarossa, S.-M. Schröder, and R. Koch, "A survey on semi-, self-and unsupervised learning for image classification," *IEEE Access*, 2021.

[18] Y. Zhao, H. Zhu, R. Liang, Q. Shen, S. Zhang, and K. Chen, "Seeing isn't believing: Towards more robust adversarial attack against real world object detectors," in *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, 2019, pp. 1989–2004.

[19] H. Wang, Y. Yu, Y. Cai, X. Chen, L. Chen, and Q. Liu, "A comparative study of state-of-the-art deep learning algorithms for vehicle detection," *IEEE Intelligent Transportation Systems Magazine*, vol. 11, no. 2, pp. 82–95, 2019.

[20] M. Carranza-García, J. Torres-Mateo, P. Lara-Benítez, and J. García-Gutiérrez, "On the performance of one-stage and two-stage object detectors in autonomous vehicles using camera data," *Remote Sensing*, vol. 13, no. 1, p. 89, 2021.

[21] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[22] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, pp. 91–99, 2015.

[23] Y. Cai, T. Luan, H. Gao, H. Wang, L. Chen, Y. Li, M. A. Sotelo, and Z. Li, "Yolov4-5d: An effective and efficient object detector for autonomous driving," *IEEE Transactions on Instrumentation and Measurement*, vol. 70, pp. 1–13, 2021.

[24] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.

[25] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[26] Ultralytics, "yolov5," [EB/OL], https://github.com/ultralytics/yolov5 Accessed on May 15, 2020.

[27] C.-Y. Wang, H.-Y. M. Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, "Cspnet: A new backbone that can enhance learning capability of cnn," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 390–391.

[28] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8759–8768.

[29] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.

[30] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu, "Traffic-sign detection and classification in the wild," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2110–2118.

[31] S.-T. Chen, C. Cornelius, J. Martin, and D. H. P. Chau, "Shapeshifter: Robust physical adversarial attack on faster r-cnn object detector," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2018, pp. 52–68.

[32] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, "Adversarial attacks on neural network policies," *arXiv preprint arXiv:1702.02284*, 2017.

[33] S. T. Jan, J. Messou, Y.-C. Lin, J.-B. Huang, and G. Wang, "Connecting the digital and physical world: Improving the robustness of adversarial attacks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 962–969.

[34] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-

image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.

[35] Y. Luo, X. Boix, G. Roig, T. Poggio, and Q. Zhao, "Foveation-based mechanisms alleviate adversarial examples," *arXiv preprint arXiv:1511.06292*, 2015.

[36] J. Lu, H. Sibai, E. Fabry, and D. Forsyth, "No need to worry about adversarial examples in object detection in autonomous vehicles," *arXiv preprint arXiv:1707.03501*, 2017.

[37] M. Xue, C. Yuan, C. He, J. Wang, and W. Liu, "Naturalae: Natural and robust physical adversarial examples for object detectors," *Journal of Information Security and Applications*, vol. 57, p. 102694, 2021.

[38] Y. Chen, X. Sun, and Y. Jin, "Communication-efficient federated deep learning with layerwise asynchronous model update and temporally weighted aggregation," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 10, pp. 4229–4238, 2019.

[39] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.

[40] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.

[41] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*, 2019.

[42] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*. Springer, 2010, pp. 177–186.

[43] M. Devi and A. Majumder, "Side-channel attack in internet of things: a survey," in *Applications of Internet of Things*. Springer, 2021, pp. 213–222.

[44] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, "Robust physical-world attacks on deep learning visual classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1625–1634.

[45] D. Song, K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, F. Tramer, A. Prakash, and T. Kohno, "Physical adversarial examples for object detectors," in *12th {USENIX} Workshop on Offensive Technologies ({WOOT} 18)*, 2018.

[46] N. Morgulis, A. Kreines, S. Mendelowitz, and Y. Weisglass, "Fooling a real car with adversarial traffic signs," *arXiv preprint arXiv:1907.00374*, 2019.

[47] G. Lovisotto, H. Turner, I. Sluganovic, M. Strohmeier, and I. Martinovic, "{SLAP}: Improving physical adversarial examples with short-lived adversarial perturbations," in *30th {USENIX} Security Symposium ({USENIX} Security 21)*, 2021.

[48] J. Li, Y. Liu, T. Chen, Z. Xiao, Z. Li, and J. Wang, "Adversarial attacks and defenses on cyber–physical systems: A survey," *IEEE Internet of Things Journal*, vol. 7, pp. 5103–5115, 2020.

[49] N. Das, M. Shanbhogue, S.-T. Chen, F. Hohman, L. Chen, M. E. Kounavis, and D. H. Chau, "Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression," *arXiv preprint arXiv:1705.02900*, 2017.

[50] C. Guo, M. Rana, M. Cisse, and L. Van Der Maaten, "Countering adversarial images using input transformations," *arXiv preprint arXiv:1711.00117*, 2017.

[51] A. Prakash, N. Moran, S. Garber, A. DiLillo, and J. Storer, "Deflecting adversarial attacks with pixel deflection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8571–8580.

[52] F. Liao, M. Liang, Y. Dong, T. Pang, X. Hu, and J. Zhu, "Defense against adversarial attacks using high-level representation guided denoiser," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1778–1787.

[53] C. Xie, Y. Wu, L. v. d. Maaten, A. L. Yuille, and K. He, "Feature denoising for improving adversarial robustness," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 501–509.

[54] X. Li and F. Li, "Adversarial examples detection in deep networks with convolutional filter statistics," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5764–5772.

[55] J. Lu, T. Issaranon, and D. Forsyth, "Safetynet: Detecting and rejecting adversarial examples robustly," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 446–454.

[56] J. H. Metzen, T. Genewein, V. Fischer, and B. Bischoff, "On detecting adversarial perturbations," *arXiv preprint arXiv:1702.04267*, 2017.

[57] W. Xu, D. Evans, and Y. Qi, "Feature squeezing: Detecting adversarial examples in deep neural networks," *arXiv preprint arXiv:1704.01155*, 2017.

[58] S. Chen, N. Carlini, and D. Wagner, "Stateful detection of black-box adversarial attacks," in *Proceedings of the 1st ACM Workshop on Security and Privacy on Artificial Intelligence*, 2020, pp. 30–39.

[59] F. Tramèr, A. Kurakin, N. Papernot, I. Goodfellow, D. Boneh, and P. McDaniel, "Ensemble adversarial training: Attacks and defenses," *arXiv preprint arXiv:1705.07204*, 2017.

[60] E. Wong, L. Rice, and J. Z. Kolter, "Fast is better than free: Revisiting adversarial training," *arXiv preprint arXiv:2001.03994*, 2020.

[61] N. Papernot, P. McDaniel, X. Wu, S. Jha, and A. Swami, "Distillation as a defense to adversarial perturbations against deep neural networks," in *2016 IEEE symposium on security and privacy (SP)*. IEEE, 2016, pp. 582–597.

[62] M. Goldblum, L. Fowl, S. Feizi, and T. Goldstein, "Adversarially robust distillation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 3996–4003.

[63] Y. Song, T. Kim, S. Nowozin, S. Ermon, and N. Kushman, "Pixeldefend: Leveraging generative models to understand and defend against adversarial examples," *arXiv preprint arXiv:1710.10766*, 2017.

[64] P. Samangouei, M. Kabkab, and R. Chellappa, "Defense-gan: Protecting classifiers against adversarial attacks using generative models," *arXiv preprint arXiv:1805.06605*, 2018.

TABLE V. DEFINITIONS OF THE NOTATIONS IN IMAGE FRAMES

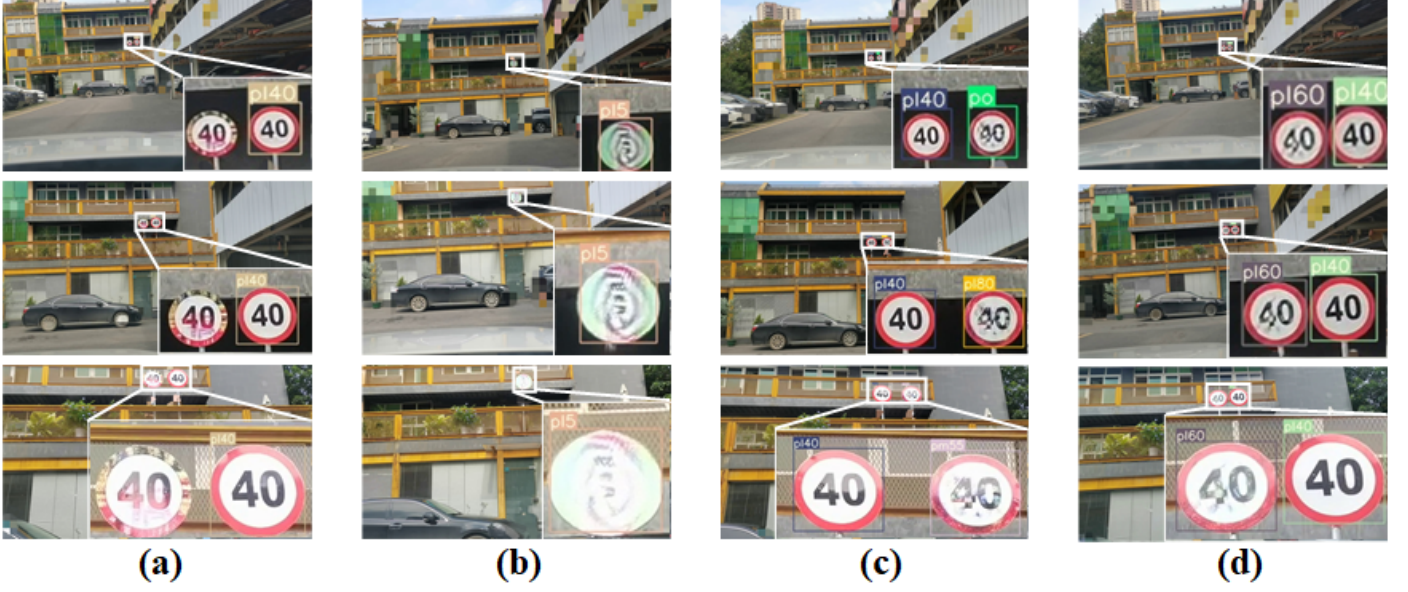| Notations in Images | pl5 | pl40 | pl60 | pl80 | pm50 | po |
|---|---|---|---|---|---|---|
| Definitions | $5km/h$ speed limit | $40km/h$ | $60km/h$t | $80km/h$ | $55ton$ weight limit | other prohibition signs |



Fig. 12. Results of four attack vectors at a height of $5.5m$. (a) HA. (b) AA. (c) NTA. (d) TA.
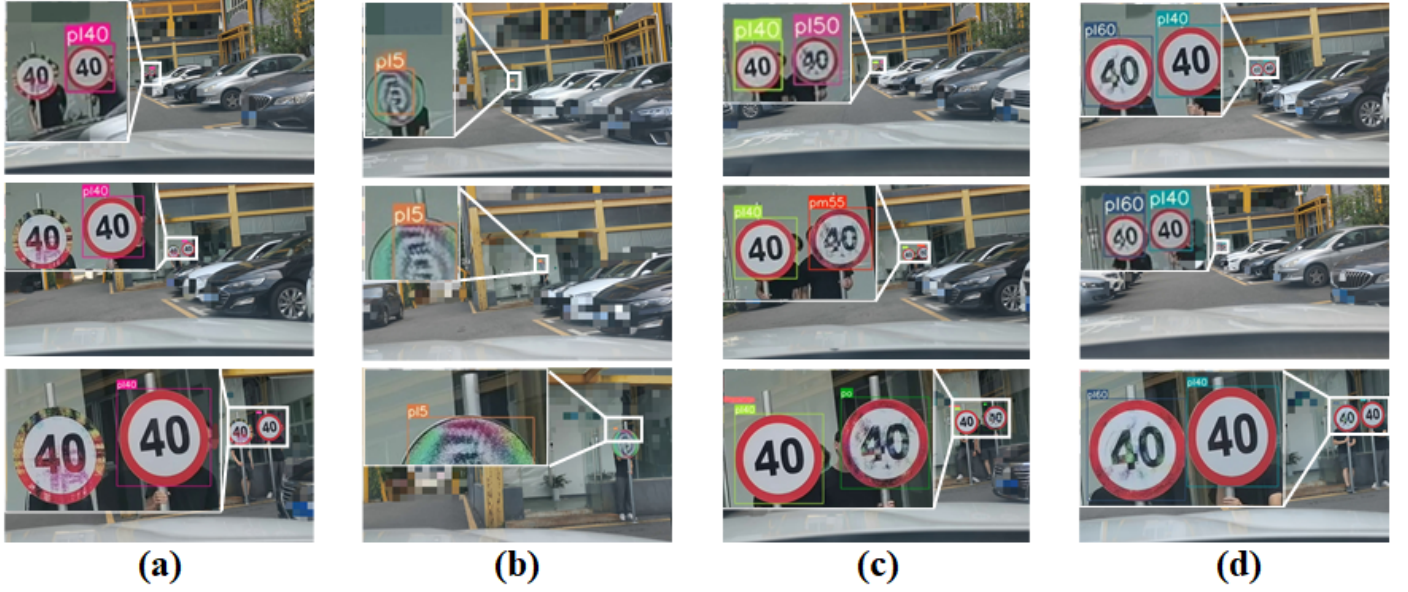


Fig. 13. Results of four attack vectors in real-road driving tests. (a) HA. (b) AA. (c) NTA. (d) TA.

Fig. 14. Results of four attack vectors at varying angles ($-60°, -45°, -30°, 0°, 30°, 45°, 60°$). (a) HA. (b) AA. (c) NTA. (d) TA.
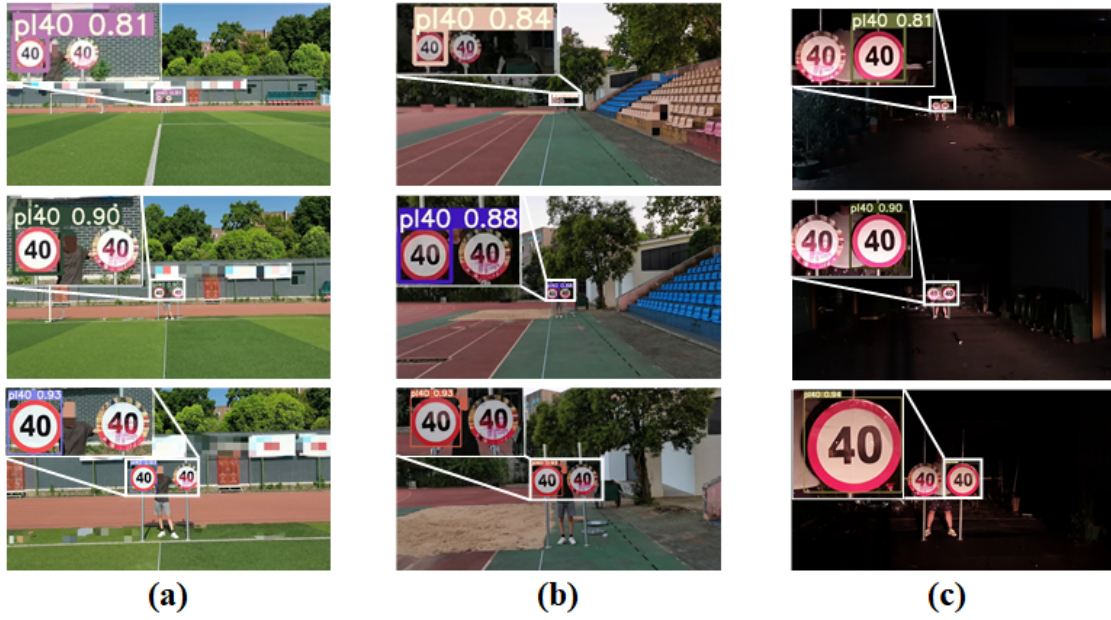


Fig. 15. Results of HA under different illuminations. (a) Sunny day. (b) Cloud day. (c) Night.