

# DARTS: Deceiving Autonomous Cars with Toxic Signs

Chawin Sitawarin\*  
Princeton University  
Princeton, NJ, USA

Arjun Nitin Bhagoji\*  
Princeton University  
Princeton, NJ, USA

Arsalan Mosenia  
Princeton University  
Princeton, NJ, USA

Mung Chiang  
Purdue University  
Princeton, NJ, USA

Prateek Mittal  
Princeton University  
Princeton, NJ, USA

## ABSTRACT

Sign recognition is an integral part of autonomous cars. Any misclassification of traffic signs can potentially lead to a multitude of disastrous consequences, ranging from a life-threatening accident to even a large-scale interruption of transportation services relying on autonomous cars. In this paper, we propose and examine security attacks against sign recognition systems for **Deceiving Autonomous caRs with Toxic Signs** (we call the proposed attacks **DARTS**). In particular, we introduce two novel methods to create these toxic signs. First, we propose **Out-of-Distribution attacks**, which expand the scope of adversarial examples by enabling the adversary to generate these starting from an arbitrary point in the image space compared to prior attacks which are restricted to existing training/test data (In-Distribution). Second, we present the **Lenticular Printing attack**, which relies on an optical phenomenon to deceive the traffic sign recognition system. We extensively evaluate the effectiveness of the proposed attacks in both *virtual and real-world settings* and consider both *white-box and black-box threat models*. Our results demonstrate that the proposed attacks are successful under both settings and threat models. We further show that Out-of-Distribution attacks can outperform In-Distribution attacks on classifiers defended using the adversarial training defense, exposing a new attack vector for these defenses.

## KEYWORDS

Adversarial examples, Autonomous cars, Lenticular printing, Security, Sign recognition, Toxic signs, Traffic signs

### ACM Reference Format:

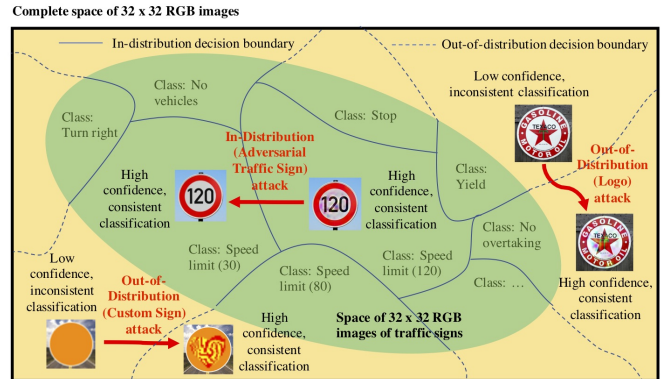
Chawin Sitawarin, Arjun Nitin Bhagoji, Arsalan Mosenia, Mung Chiang, and Prateek Mittal. . DARTS: Deceiving Autonomous Cars with Toxic Signs. In *Proceedings of* . ACM, New York, NY, USA, 18 pages.

## 1 INTRODUCTION

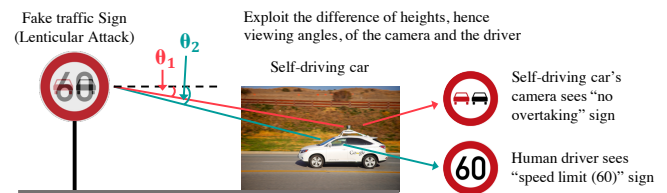
The rapid technological and scientific advancements in artificial intelligence (AI) and machine learning (ML) have led to their deployment in ubiquitous, pervasive systems and applications, such

\*Both authors contributed equally to the paper

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).



(a) Illustration of Out-of-Distribution evasion attacks on a traffic sign recognition system trained with traffic sign images. Out-of-Distribution attacks enable the adversary to start from anywhere in the space of images and do not restrict her to the training/test data.

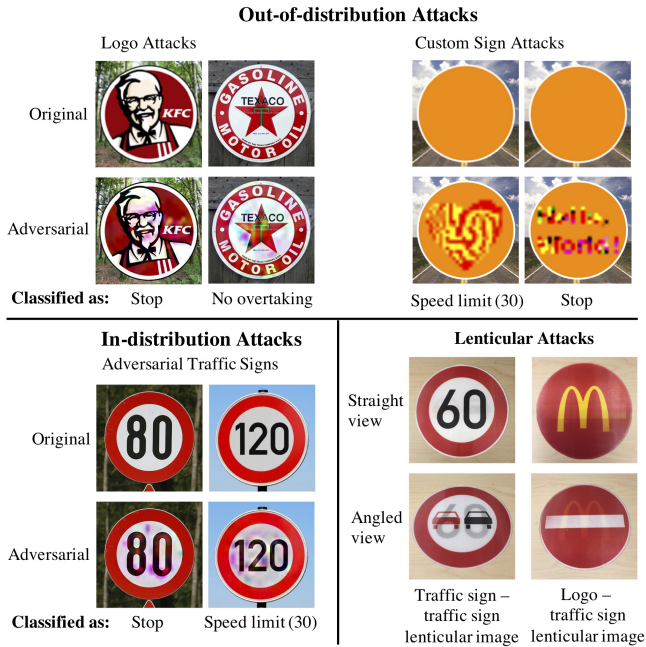


(b) Principle behind an attack based on lenticular printing.

Figure 1: New attack vectors

as authentication systems [1, 2], health-care applications [3, 4], and several vehicular services [5–8]. The ubiquity of ML provides adversaries with both opportunities and incentives to develop strategic approaches to fool learning systems and achieve their malicious goals [9, 10]. A number of powerful attacks on the test phase of ML systems used for classification have been developed over the past few years, including attacks on Support Vector Machines [11, 12] and deep neural networks [12–17]. These attacks have also been shown to work in black-box settings [18–23]. These attacks work by adding carefully-crafted perturbations to benign examples to generate adversarial examples. In the case of image data, these perturbations are typically imperceptible to humans. *While these attacks are interesting from a theoretical perspective and expose gaps in our understanding of the working of neural networks, their practical importance has remained unclear in real-world application domains.*

A few recent research studies have attempted to provide realistic attacks on ML classifiers, which interact with the physical world



**Figure 2: Toxic signs for our traffic sign recognition pipeline generated using the In-Distribution attack, the Out-of-Distribution attack (Logo and Custom Sign), and the Lenticular Printing attack. The adversarial examples are classified as the desired target traffic sign with high confidence under a variety of physical conditions when printed out. The Lenticular Printing attack samples flip the displayed sign depending on the viewing angle, simulating the view of the human driver and the camera in the autonomous car.**

[24, 25]. Kurakin et al. [24] print out virtual adversarial examples, take their pictures, and then pass them through the original classifier. Sharif et al. [25] attack face recognition systems by allowing subjects to wear glasses with embedded adversarial perturbations that can deceive the system. However, both these attacks were carried out in *controlled laboratory setting* where variations in brightness, viewing angles, distances, image re-sizing etc. were not taken into account, imposing limitations on their effectiveness in real-world scenarios. Further, they only focus on *creating adversarial examples from training/testing data* (e.g., *image of different faces in the face recognition system*) to attack the underlying systems. We refer to these as In-Distribution attacks. Athalye et al. [26] introduced the Expectation over Transformations (EOT) method to generate physically robust adversarial examples and concurrent work by Evtimov et al. [27] used this method to attack traffic sign recognition systems. A detailed comparison of our work with these is at the end of the introduction.

In this paper, we focus on physically-realizable attacks against sign recognition system utilized in autonomous cars, one of the most important upcoming application of ML [5, 6, 28, 29]. We introduce two novel types of attacks, namely, *Out-of-Distribution* and *Lenticular Printing attacks* that deceive the sign recognition system, leading to potentially life-threatening consequences. Our novel

attacks shed lights on *how domain-specific characteristics of an application domain, in particular, autonomous cars, can be exploited by an attacker to target the underlying ML classifier.*

Our key contributions can be summarized as follows:

**1. Introducing new attack vectors:** We introduce new methods to create toxic signs that look benign to human observers, at the same time, deceive the sign recognition mechanism. Toxic signs cause misclassifications, potentially leading to serious real-world consequences including road accidents and widespread traffic confusion. We significantly extend the scope of attacks on real-world ML systems in two ways (Figure 1). First, we propose **Out-of-Distribution attack**, which enables the adversary to start from an arbitrary point in the image space to generate adversarial examples, as opposed to prior attacks that are restricted to samples drawn from the training/testing distribution. The proposed attack is motivated by the key insight that autonomous cars are moving through a *complex environment consisting of many objects, which can potentially used by an attacker to create adversarial examples.* We previously provided a high-level description of Out-of-Distribution attack in an Extended Abstract [30]. In this paper, we provide an in-depth explanation of this attack and thoroughly examine its effectiveness in various experimental scenarios. Second, we present **Lenticular Printing attack**, which relies on an optical phenomenon to fool the sign recognition system. Figure 2 demonstrates a few toxic signs created by the two above-mentioned attacks.

**2. Extensive experimental analysis:** We evaluate the proposed attacks in both *virtual and real-world settings* over various sets of parameters. We consider both the *white-box threat model* (i.e., *the adversary has access to the details of the traffic sign recognition system*) and the *black-box one* (i.e., *such access is not present*). We demonstrate that adversarial examples (created from either arbitrary points in the image space or traffic signs) can deceive the traffic sign recognition system with high confidence. Further, we show the attacker can achieve significant attack success rates even in black-box settings. To provide a thorough analysis of the proposed attacks, we also conduct real-world drive-by tests, where a vehicle-mounted camera continuously captures image from the surroundings and offers its data to the sign recognition system (Figure 8). We achieve attack success rates in excess of 90% in the real-world setting with both Out-of-Distribution and In-Distribution attacks.

**3. Studying the effect of Out-of-Distribution attacks on state-of-the-art defenses:** We discuss the limitations of adversarial training based defenses [14] in mitigating the proposed attacks. We show Out-of-Distribution attacks, in which the initial image does not come from the underlying training/testing distribution, outperform In-Distribution attacks on adversarial training [14] based defenses in which adversarial examples are created based on the initial training dataset and are considered in the training phase. Known ML-based defenses are intrinsically prone to the Lenticular Printing attack which has been developed with respect to the physical characteristics of the application domain.

**Comparison with Athalye et al. [26] and Evtimov et al. [27]:** Athalye et al. [26] introduced the Expectation over Transformations (EOT) method to generate physically robust adversarial examples. They used the method to generate 3D printed adversarial examples which remain adversarial under a range of conditions. These samples are evaluated on classifiers trained on the Imagenet dataset.

In *concurrent* work, Evtimov et al. [27] used the EOT method to generate physically robust adversarial ‘Stop’ and ‘Right Turn’ signs in a white-box, In-Distribution setting. They manually crop out the portion of the video frame corresponding to the adversarial example while we use an automated detection pipeline. We further expand the space of adversarial examples in both virtual and real-world settings by introducing Out-of-Distribution attacks. We also evaluate the effectiveness of transferability-based black-box attacks for both In-Distribution and Out-of-Distribution settings, as well as defenses using adversarial training. We also introduce the completely new attack vector of Lenticular Printing attacks based on optical phenomena which exist outside the ambit of adversarial examples.

Our proof-of-concept attacks shed light on fundamental security challenges associated with the use of sign recognition techniques in autonomous cars, paving the way for further investigation of overlooked security challenges in this domain.

## 2 BACKGROUND

In this section, we present relevant background on machine learning systems, the traffic sign recognition pipeline and the adversarial examples and threat models we consider.

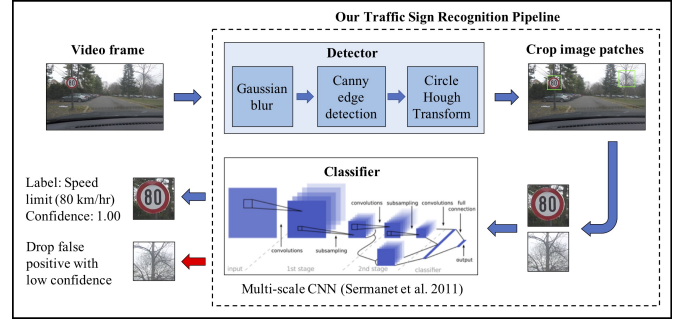
### 2.1 Supervised machine learning systems

Machine learning systems typically have two phases, a training phase and a test phase [31]. A classifier  $f(\cdot; \theta) : \mathcal{X} \rightarrow \mathcal{Y}$  is a function mapping from the domain  $\mathcal{X}$  to the set of classification outputs  $\mathcal{Y}$ . The number of possible classification outputs is then  $|\mathcal{Y}|$ .  $\theta$  is the set of parameters associated with a classifier.  $\ell_f(\mathbf{x}, y)$  is used to represent the loss function for the classifier  $f$  with respect to inputs  $\mathbf{x} \in \mathcal{X}$  and labels  $y \in \mathcal{Y}$ . The classifier is trained by minimizing the loss function over  $n$  samples  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  drawn from a distribution  $\mathcal{D}$  over the domain  $\mathcal{X}$ . In the particular case of *traffic sign recognition*,  $\mathcal{D}$  is a distribution over images of traffic signs.

Since *deep neural networks* (DNNs) [32] achieve very high classification accuracy in a variety of image classification settings [33, 34], they are being used for the computer vision systems of autonomous cars [28, 29]. We thus focus on attacks on DNNs in this paper and define some notation specifically for neural networks. The outputs of the penultimate layer of a neural network  $f$ , representing the output of the network computed sequentially over all preceding layers, are known as the logits. We represent the logits as a vector  $\phi^f(\mathbf{x}) \in \mathbb{R}^{|\mathcal{Y}|}$ . The final layer of a neural network  $f$  used for classification is usually a softmax layer. The loss functions we use are the standard *cross-entropy loss* [32] and the *logit loss* [15].

### 2.2 Traffic Sign recognition pipeline

Our traffic sign recognition pipeline (Figure 3) consists of two stages: detection and classification. We utilize a commonly used recognition pipeline based on the Hough transform [35–37]. The shape-based detector uses the circle Hough transform [38] to identify the regions of a video frame or still image that contain a circular traffic sign. Before using Hough transform, we smooth a video frame with a Gaussian filter and then extract only the edges with Canny edge detection [39]. Triangular signs can be detected by a similar method



**Figure 3: Sign recognition pipeline for real-world evaluation. The pipeline consists of an initial detection phase followed by a multi-scale CNN as a classifier. In the virtual setting, a video frame is replaced by a still image.**

[36]. The detected image patch is cropped and re-sized to the input size of the classifier before it is passed on to the neural network classifier trained on a traffic sign dataset. The classifier outputs confidence scores for all output classes to determine whether the input is a traffic sign and assign its label. The label with the highest confidence is chosen as the final output only if its confidence is above a certain threshold. Images classified with a low confidence score are discarded as false positives for detection.

### 2.3 Adversarial examples and threat models

Our focus is on attacks during the test phase, which are typically known as *evasion attacks*. These have been demonstrated in the virtual setting for a number of classifiers [12, 14, 15, 17, 40, 41]. These attacks aim to modify benign examples  $\mathbf{x} \in \mathcal{X}$  by adding a perturbation to them such that the modified examples  $\tilde{\mathbf{x}}$  are *adversarial*, i.e. they are misclassified by the ML system in a targeted class (targeted attack), or any class other than the ground truth class (untargeted attack). In the case of attacks on the computer vision systems of autonomous cars, we focus entirely on *targeted attacks* since these are more realistic from an attacker’s perspective. To generate a *targeted* adversarial sample  $\tilde{\mathbf{x}}$  starting from a benign sample  $\mathbf{x}$  for a classifier  $f$ , the following optimization problem [15] leads to state-of-the-art attack success rates in the virtual setting:

$$\begin{aligned} \min \quad & d(\tilde{\mathbf{x}}, \mathbf{x}) + \lambda \ell_f(\tilde{\mathbf{x}}, T), \\ \text{s.t.} \quad & \tilde{\mathbf{x}} \in \mathcal{H}. \end{aligned} \quad (1)$$

We use this attack as a baseline. Here  $d$  is an appropriate distance metric for inputs from the input domain  $\mathcal{X}$  (usually an  $L_p$  norm),  $T$  is the target class and  $\mathcal{H}$  is the constraint on the input space.  $\lambda$  controls the trade-off between minimizing the distance to the adversarial example and minimizing the loss with respect to the target. In essence, the optimization problem above tries to find the closest  $\tilde{\mathbf{x}}$  to  $\mathbf{x}$  such that the loss of the classifier at  $\tilde{\mathbf{x}}$  with respect to the target  $T$  is minimized. For a neural network,  $\ell_f(\cdot, \cdot)$  is typically highly non-convex, so heuristic optimizers based on stochastic gradient descent have to be used to find local minima [15].

For traffic sign recognition systems, the method described above produces adversarial examples which do not work well under conditions encountered in the real world such as variation in brightness,



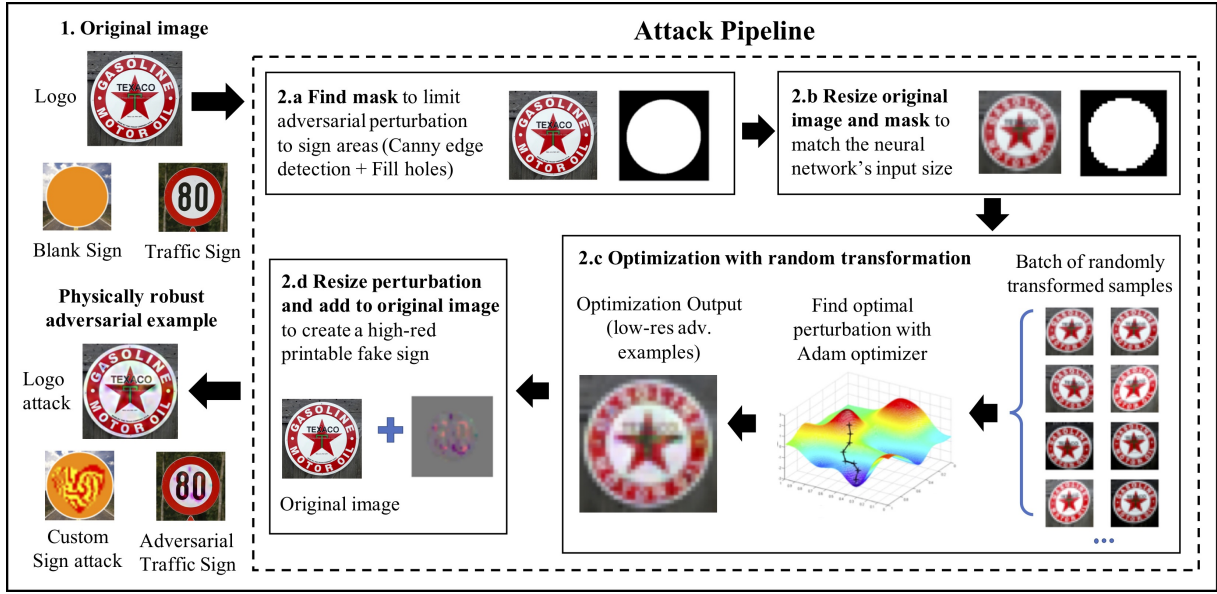


Figure 4: Overview of the Attack pipeline. This diagram provides an overview of the process by which adversarial examples are generated for both the Out-of-Distribution attack and In-Distribution attacks.

viewing angles, image re-sizing etc.. That is, if the examples generated using this method were directly used in Step 3 (mount sign and drive-by test) of Figure 4, they would not perform well. In light of this, we incorporate these conditions directly into a modified optimization problem used in Section 3 to generate physically robust adversarial examples. Similar approaches to this problem have been taken in Athalye et al. [26] as well as in concurrent work by Evtimov et al. [27].

**2.3.1 Threat models.** We consider two commonly used threat models for the generation of adversarial examples against deep neural networks: *white-box* and *black-box*. The black-box setting also applies to the Lenticular Printing attack.

**White-box:** In the white-box setting, we assume that the adversary has complete access to the target model  $f$  including its architecture and weights. We briefly justify the powerful attacker considered in the white-box setting. Consider an attacker who wishes to cause an autonomous car to detect and misclassify a sign in its environment. It is conceivable that the attacker can purchase or rent a vehicle of the kind that is to be attacked, and ‘reverse engineer’ the classifier by querying it on an appropriate dataset to train a surrogate classifier that closely mimics the target classifier [42]. Further, direct query-based attacks can be as powerful as white-box attacks [22, 23]. Thus, the white-box setting is an important one to consider.

**Black-box (no query access):** In this setting, the adversary does not have direct access to the target model. We do not even assume query access as in previous work [21–23]. In this setting then, black-box attacks rely on the phenomenon of transferability [13, 18–20], where adversarial examples generated for a model trained locally by the attacker, remain adversarial for a different, target model. The Lenticular Printing attack naturally operates in the black-box setting since it relies on an optical phenomenon and not on the internal structure of the classifier.

### 3 ATTACKS: ADVERSARIAL EXAMPLES

In this section, we present new methods to generate physically robust adversarial examples. The main aim for these examples is to be *classified consistently as the desired target class under a variety of real-world conditions* such as variations in brightness, viewing angle, distance and image re-sizing. In particular, we introduce *Out-of-Distribution attacks*, which modify arbitrary, innocuous signs such that they are detected and classified as potentially dangerous traffic signs in the attacker’s desired class. This attack *greatly enlarges the possible space of adversarial examples* since the adversary now has the *ability to start from any point in the space of images to generate an adversarial example*. We also examine the In-Distribution attack which generates adversarial examples starting from existing traffic signs. We then describe the pipeline we use to ensure these generated adversarial examples are classified consistently even under randomized input transformations which we use to model real-world conditions that may be encountered.

#### 3.1 Attack overview

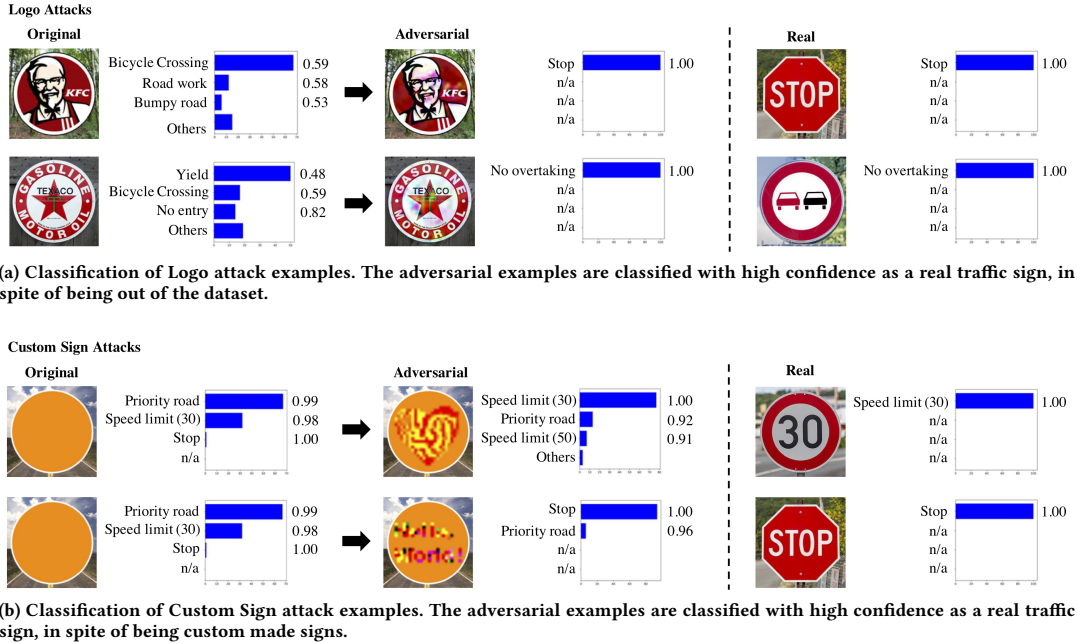
We first provide an overview of our attack pipeline from Figure 4 and then describe each of its components in detail in the following subsections. Our pipeline has three steps:

**Step 1.** Obtain the original image  $x$  and choose target class  $T$  that the adversarial example  $\tilde{x}$  should be classified as. In our attacks,  $x$  can either be an *in-distribution image* of a traffic sign or an *out-of-distribution image* of an ad sign, logo, graffiti etc.

**Step 2.** Generate the digital version of the physically robust adversarial example as follows:

1. Generate mask  $M$  for the original image (A mask is needed to ensure that the adversary’s perturbation budget is not utilized in adding perturbations to the background.)





**Figure 5: Frequency of the top-3 labels the given images are classified as under 100 different randomized transformations. The numbers beside each bar on the bar charts provides the average confidence of the top classification outcomes over the 100 randomized transformations.**

2. Re-size both the original image and the mask to the input size of the target classifier<sup>1</sup>.

3. Run the physically robust adversarial example optimization problem from Equation 2 to obtain the perturbation  $\delta$ . The optimization problem includes randomized brightness, perspective (includes rotations and shearing) and size variations while generating the adversarial example to ensure it remains adversarial under real-world conditions.

4. Re-size the output perturbation  $\delta$  and add it to the original image.

**Step 3.** Print and test the generated adversarial example  $\tilde{x}$  for robustness in real-world conditions.

### 3.2 Step 1: Choosing the input

Now, we describe two different attack modes in which an adversary can operate to generate adversarial examples: (1) the Out-of-Distribution attack, where the adversary is free to generate adversarial examples from any innocuous elements in the environment such as advertisement signs, drawings, graffiti etc. and (2) the In-Distribution attack, where the adversary modifies existing traffic signs to make them adversarial. The second setting is similar to the one considered in most previous work on generating adversarial examples [14, 15].

**3.2.1 Out-of-Distribution attacks.** We propose a novel attack based on the concept of adversarial examples by exploiting the fact that any image in the input domain of the classifier can be made adversarial by adding an imperceptible perturbation generated by

<sup>1</sup>This is  $32 \times 32$  pixels for the classifiers we use, described in Table 1

an optimization problem such as Equation 2 which also ensures physical robustness. Since the classifier is only trained on images of traffic signs, it can only be expected to reliably classify other traffic signs, which are effectively, *in distribution* for that classifier. However, the fact that it provides a classification outcome for *any input image*, represents a security risk which we exploit in our attacks. In particular, we start with an *out-of-distribution image* (not a traffic sign) and generate a targeted adversarial example from it. Here, we demonstrate two possible instantiations of the Out-of-Distribution attack.

**Logo attacks:** In this attack, images of commonly found logos are modified such that they are detected and classified with high confidence as traffic signs (Figure 5a). Since these logos are omnipresent, they allow an adversary to carry out attacks in real-world settings such as city streets. In this scenario, the attack pipeline from Section 3.1 is used and the adversarial perturbation is constrained to be as small as possible while still being effective under transformations.

**Custom Sign attacks:** In this attack, the adversary creates a custom sign that is adversarial starting from a blank sign (Figure 5b). Any mask corresponding to graffiti, text etc. on blank signs can lead to the embedding of adversarial traffic signs in inconspicuous, graffiti-like objects in the environment. This allows the adversary to create adversarial examples in any setting by using a mask to create images or text that are appropriate for the surroundings. In this attack, the original sign is a solid color circular sign and the norm of the perturbation is not penalized by the optimization problem but only constrained by the shape of the mask. This allows the adversary to draw almost any desired shape on a blank sign

and the optimization will try to fill out the shape with colors that make the sign classified as the target class.

Under ordinary conditions when no adversarial examples are present, the false recognition of objects which are not traffic signs does not affect the traffic sign recognition system since (i) the confidence scores corresponding to the predicted labels of these objects are usually low; (ii) these circular objects are not consistently classified as a certain sign. The predicted label changes randomly as the background and viewing angle varies across multiple frames in the video.

Therefore, a traffic sign recognition system, including ours, can choose to treat any detection with these two properties as an erroneous detection by setting the confidence threshold close to 1 and/or ignoring objects that are inconsistently classified. On the other hand, the generated adversarial examples are *classified consistently as target traffic signs with high confidence under varying physical conditions* which we demonstrate experimentally below.

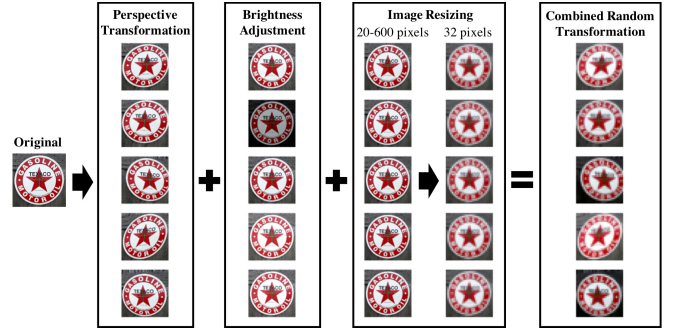
**Confirming hypothesis about classifier confidence:** To confirm our earlier hypotheses with regard to the confidence of classification for Out-of-Distribution images and their adversarial counterparts, we display some adversarial examples, their classification and confidence in Figure 5. Experimental details are in Section 6.

We apply randomized transformations to create 100 images for each of the Logo signs. The "Original" column of Figure 5a shows that the logo signs are classified as different classes depending on the transformation and with low confidence on average. As a comparison, the "Real" column of Figure 5 shows that real traffic signs are consistently classified as the correct label with probability very close to 1. The "Adversarial" column demonstrates that the generated adversarial examples are *consistently classified with high confidence* in the desired target class. In a sense, the adversarial logo signs and the real traffic signs are equivalent from the perspective of the classifier.

Similarly, the Custom Sign attack samples are mostly classified as the target class with high confidence under 100 different randomized transformations. Note that even though the original solid color signs shown in the "Original" column of Figure 5b are classified with high confidence, their classified labels are arbitrary and not under the control of the adversary. They are highly dependent on the signs' orientation and brightness and cannot be used as reliable adversarial examples. For example, a slight change in the camera angle may yield different labels which is an undesirable effect from the adversary's perspective.

**Note:** Out-of-Distribution attacks can be carried out against any classifier in any application domain. We choose to focus on traffic sign recognition systems since they provide a compelling setting where such attacks are plausible and effective. Another application domain of interest is the computer vision systems of Augmented Reality systems. In certain other settings, such as fooling content moderation systems [22], Out-of-Distribution attacks may be of limited interest.

**3.2.2 In-Distribution attack.** In this attack, images of *traffic signs* are modified using imperceptible perturbations such that they are classified as a different traffic sign. This attack is similar to most attacks carried out in most previous work in both the virtual [12–15,



**Figure 6: Transformations used during both training and evaluation to simulate real-world conditions. The final column of images displays the type of images that are included in the sum in Eq. (2).**

[17] and physical [24–27] settings. We include it here to demonstrate that our pipeline works in this setting as well.

### 3.3 Step 2: Robust adversarial example generation

In this section, we describe our methodology for generating robust, physically realizable adversarial examples starting from either of the inputs described in the previous section.

**3.3.1 Optimization problem.** Our adversarial example generation method involves heuristically solving a non-convex optimization problem using the Adam optimizer [43] to find the optimal perturbation  $\delta$ . The problem set-up is adapted from the general concept of expectation over transformations [26]. An updating gradient is averaged over a batch of randomly transformed versions of the original image [27, 44]. The robust adversarial example can be written as a solution to the minimization problem given below for any input  $\mathbf{x}$ :

$$\min_{\delta \in \mathbb{R}^n} c \cdot \frac{1}{B} \sum_{i=1}^B [F(\tau_i(\mathbf{x} + M \cdot \delta))] + \max(\|\delta\|_p, L) \quad (2)$$

where  $F(\mathbf{x}) = \max(\max_{j \neq T} \{\phi(\mathbf{x})_j\} - \phi(\mathbf{x})_T, -K)$  is the logit loss [15] and  $\phi(\mathbf{x})_j$  is the  $j^{\text{th}}$  logit of the target network. The constant  $K$  determines the desired objective value controls the *confidence score* of the adversarial example.  $M$  is a mask or a matrix of 0s and 1s with the same width and height as the input image which is multiplied element-wise with the perturbation ( $M \cdot \delta$ ) to constrain the feasible region to the sign area.  $\tau_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$  is a transformation function mapping within the image space ( $n = 32 \times 32 \times 3$ ). The transformations  $\tau_i$  are chosen to be differentiable so that gradients can be propagated back to the original image.

The objective value is penalized by a  $p$ -norm of the perturbation  $\delta$  in order to keep the adversarial perturbation imperceptible by humans, and the constant  $c$  is adjusted accordingly to balance between the loss and penalty terms. We introduce an additional constant  $L$  to explicitly encourage the norm of the perturbation to be at least  $L$  since *overly small perturbations can disappear or be rendered ineffective in the process of printing and video capturing*. In practice for the Custom Sign attack, the same optimization problem is used by setting  $c$  and  $L$  to some large numbers so that the

optimization will focus on minimizing the loss function without penalizing the norm of the perturbation.

**Summary:** The optimization problem tries to minimize the average of the objective function with respect to a perturbation  $\delta$  restricted to an area determined by a mask  $M$  evaluated on a batch of  $B$  images which are generated from a set of transformation functions  $\tau_1, \dots, \tau_B$ .

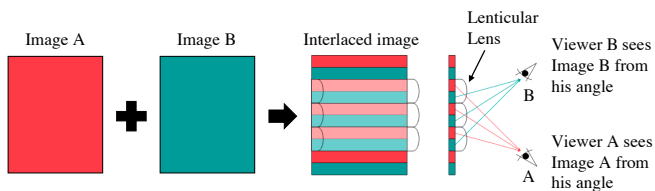
**3.3.2 Image selection and re-sizing.** We use a single RGB image, re-sized to the input size of the classifier. We use either bilinear interpolation or nearest neighbours for re-sizing. This input image is then transformed by the set of differentiable and predefined transformations with randomized parameters. The attacker has a low cost with regards to constructing an adversarial example since only a single image is needed to construct robust adversarial examples. In the case of In-Distribution attacks, the adversary only needs one image of the same class as the original sign which can be easily found on the internet. The image also does not need to be a real photograph; as we show in our Out-of-Distribution attacks, any logo or schematic drawing can be used as well.

**3.3.3 Mask generation.** A mask is a mapping from each pixel of the image to 0 or 1. Background pixels are mapped to 0 while pixels from the sign area are mapped to 1. To create a *mask*, we employ a simple image segmentation algorithm<sup>2</sup> using Canny edge detection [46] to outline the boundary of the sign and binary dilation to fill in the hole. This algorithm works well on an image of a sign of any shape, given that it is of a high resolution.

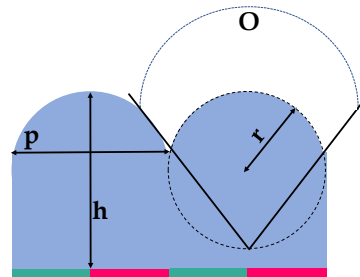
**3.3.4 Transformations.** Here, we describe the set of image transformations used in the optimization which are chosen to virtually generate a batch of input images resembling the varying real-world conditions under which the generated physical adversarial example is expected to be effective. We included three different differentiable transformations in our experiments: (1) perspective transform, (2) brightness adjustment and (3) re-sizing, each of which is randomized for a particular image while generating an adversarial example as well as while evaluating it (see metrics used in Section 5.2.1). Examples of transformations used for both the optimization and evaluation phases are shown in Figure 6.

**Perspective transformation:** Orientations of a sign that appears on an image vary between camera angles, distances and relative heights. These orientations can be captured by perspective transformation, an image transformation that maps each of the four corners of a 2D image to a different point in the 2D space. Any perspective transformation can be characterized by eight degrees of freedom which can be represented by a  $3 \times 3$  *projective transform matrix* with the last entry being set to 1. Other common image transformations such as rotation, shearing or scaling are subsets of the perspective transformation.

**Brightness adjustment:** Naturally, cars drive any time of the day so an adversary might need the fake sign to be able to fool the system under different lighting conditions. The simplest way to simulate settings with different amounts of light or brightness is to add a constant value to every pixel and every channel (R, G, B) of an image and then clip the value to stay in the allowed range of



(a) Illustration of the process of generating a lenticular image and its angle-dependent appearance.



(b) Several parameters ( $p$ ,  $r$ , and  $h$ ) determine the full angle of observation ( $O$ )

**Figure 7: Concepts underlying the working of the Lenticular Printing attack.**

[0, 1]. This transformation consisting of an addition and clipping is differentiable almost everywhere.

**Image re-sizing:** Due to varying distances between a camera and a sign, an image of the sign detected and cropped out of the video frame could have different sizes (equivalent to number of pixels). Further, after generation, adversarial perturbations are up-sampled to match the original image size. Subsequently, when the adversarial sign is detected on the camera, it is cropped out and down-sampled back to  $32 \times 32$  pixels for the classifier. This re-sampling process could blur or introduce artifacts to the image, diminishing the intended effect of adversarial perturbations. We use a *nearest neighbours* re-sizing technique for the mask and *bilinear interpolation* for all other images.

The randomized transformation  $\tau$  used during both the generation and evaluation of adversarial examples is a composition of the three transformations mentioned above. The degree of randomness of each transformation is controlled by a set of tunable parameters to keep them in a range that provides realistic simulations.

**Parameter choices:** Our optimization program introduces a number of parameters ( $c, K, L, T, \theta$ , step size, choices of norm) that can be fine-tuned to control the trade-off between robustness of adversarial examples and their visibility to humans. A detailed analysis of the procedure we followed to optimize these parameters in tandem to achieve high attack success rates in practice is in Appendix E.

## 4 ATTACKS: LENTICULAR PRINTING

In this section, we describe a novel attack on traffic sign recognition systems that we call the Lenticular Printing attack.

**Key idea:** Lenticular Printing attack is motivated by the key insight that the human driver and the vehicle-mounted camera observe

<sup>2</sup>Code adapted from [45]



Component	Name	Description	Usage
Datasets	GTSRB [47]	51,839 $32 \times 32$ RGB images of 43 types of traffic signs (39,209 for training, 12,630 for testing) (Figure 10)	Training and validating both classifiers
	GTSDB [48]	900 $1360 \times 800$ RGB images of traffic signs in real-world conditions	Evaluating detector+classifier pipeline
	Auxiliary traffic data	22 high-resolution RGB images of traffic signs (Figure 11)	Generating real-world examples for the In-Distribution attack
	Logo	7 high-resolution RGB images of popular logos (Figure 12)	Generating real-world examples for the Out-of-Distribution attack
	Custom Sign	3 blank circular signs (orange, white and blue) combined with 6 custom masks to constrain perturbations (Figure 13)	Generating real-world examples for the Out-of-Distribution attack
Classifiers	Multi-scale CNN [49]	Multi-scale CNN with 3 convolutional layers and 2 fully connected layers <b>98.50%</b> accuracy on GTSRB, <b>77.1% mAP</b> on GTSDB, 100% accuracy on auxiliary traffic data	Main target network (virtual and real-world)
	Standard CNN	Conventional CNN with 4 convolutional layers and 3 fully connected layers <b>98.66%</b> accuracy on GTSRB, <b>81.5% mAP</b> on GTSDB, 100% accuracy on auxiliary traffic data	Evaluating transferability-based black-box attacks

**Table 1: Summary of classifiers and datasets used for experimental evaluation.**

the environment from two different observation angles (Figure 1). To deceive the sign recognition system utilized in autonomous cars, we create special traffic signs that appear differently from different observation angles (Figure 2). In particular, we exploit a multi-step process, known as lenticular printing, that involves creating a special image from at least two existing images, and combining it with an array of magnifying lenses (Figure 7a). Lenticular printing relies on an optical phenomenon and has been traditionally used in photography and visual arts to generate 2-D images that can offer an illusion of depth and be changed as the image is viewed from different angles.

**Choosing an appropriate lens array:** Here, we briefly describe how an attacker can use this technique to create malicious traffic signs. To maximize the possibility of observing two different images from two different angles, the attacker should interlace images so that the width of each row of an image (red and green rows in Figure 7a) is equal to half of the width of each lens in the lens array, i.e.,  $Width_{rows} = p/2$ , where  $p$  is shown in 7b. For a successful attack, the attacker must ensure that the difference between the observation angle of the driver and the observation angle of the camera always remains between  $O/4$  and  $O/2$  (i.e.,  $O/4 \leq |\theta_1 - \theta_2| \leq O/2$ ), where  $\theta_1 - \theta_2$  are shown in Figure 1b and  $O$  is the full angle of observation of a lenticular lens (Figure 7b) and can be determined as follows:

$$O = 2(\sin^{-1}(\frac{p}{2r}) - \sin^{-1}(\frac{n * \sin(\sin^{-1}(\frac{p}{2r}) - \tan^{-1}(\frac{p}{h}))}{n_a})), \quad (3)$$

where  $n_a$  is 1.003 (the index of refraction of air),  $n$  is the index of refraction of the lens and all other parameters are shown in Figure 7b. In our experiments, we use commonly-available lens with a  $49 \sim 50$ -degree full angle of observation (i.e.,  $O \approx 50 \Rightarrow 12.5 < |\theta_1 - \theta_2| < 25$  should be maintained during driving).

**Creating malicious signs:** Industrial quality lenticular printing involves specialized machinery and trained personnel to operate it. Nevertheless, a simple lenticular image can also be produced without the need for specialized equipment by using either an inkjet or a laser color printer and a lenticular lens. We have created lenticular images using a two-step procedure: (i) we obtain two images of the same dimensions, one of which is a standard traffic sign (the adversary’s desired target) and the other one can be a logo, a custom sign, or another traffic sign. These just have to be of the same dimensions so that they can be interlaced with the desired traffic sign; (ii) we print the interlaced image on photo-quality paper and stick it on the back of a commercially available lenticular lens. We have used a free software called “SuperFlip” available online to interlace the chosen images [50].

## 5 EXPERIMENTAL GOALS AND SETUP

In this section, we describe our experimental goals and setup, which motivate the results in the subsequent sections.

### 5.1 Experimental goals

It is well-known that In-Distribution adversarial examples are extremely effective in the virtual setting. With our experiments, we show that Out-of-Distribution attacks are effective in both virtual and real-world settings. In particular, we seek to demonstrate with an appropriate pipeline for the generation, evaluation and selection of adversarial examples, an adversary can achieve high real-world attack success rates. In particular, we answer the following questions in our experiments:

- (1) How effective are Out-of-Distribution attacks starting from high-resolution real-world images on standard classifiers? Answer: Section 6
- (2) How does an adversary generate and evaluate adversarial examples for a real-world setting? Answer: Section 6.3
- (3) Can Out-of-Distribution attacks break existing state-of-the-art defenses for neural networks? Answer: Section 6.4
- (4) Are transferability-based black-box attacks effective in real-world settings? Answer: Section 7
- (5) Can the Lenticular Printing attack fool classifiers? Answer: Section 7.1

Overall, we demonstrate that with our pipeline, it is a low overhead procedure for an adversary to test both In-Distribution and Out-of-Distribution adversarial examples and subsequently use these to carry out effective real-world attacks.

### 5.2 Experimental setup

For our experiments we used a GPU cluster of 8 NVIDIA Tesla P100 GPUs with 16GB of memory each, running with no GPU parallelization. A summary of the datasets and classifiers used is given in Table 1. The overall pipeline was described earlier in Figure 3. Dataset details are in Appendix A, detector details in Appendix B and classifier details in Appendix C.

**5.2.1 Metrics.** Our experiments only consider *targeted attacks*. **Virtual attack success (VAS):** This is the standard evaluation metric used in previous work [15, 20, 22] where the success rate is measured as the fraction of adversarial examples that are classified as their target.

**Simulated physical attack success (SPAS):** Manually printing out and evaluating the effectiveness of adversarial examples in a real-world setting is prohibitively expensive. In light of this, we

Attacks	Virtual attack success (VAS)	Simulated physical attack success (SPAS)	Avg. norm ( $L_1$ )	Avg. confidence
In-Distribution (auxiliary traffic data)	54.34 %	36.65 %	37.71	0.9721
Out-of-Distribution (Logo)	85.71%	65.07%	34.89	0.9753
Out-of-Distribution (Custom Sign)	29.44%	18.72%	N.A.	0.9508

**Table 2: White-box attack success rates for In-Distribution and Out-of-Distribution attacks on the Multi-scale CNN (virtual setting).**

propose a method to evaluate how physically robust adversarial examples are likely to be by simulating varying conditions in which the images of them might be taken. The physical conditions are virtually simulated by a combination of the randomized image transformations described earlier (Section 3.3.4). *10 randomized, composite transformations* (brightness, perspective and re-sizing) are applied to each adversarial and original sample, and the transformed images are directly fed into the target classifier to determine the predicted labels. The *simulated physical attack success* (SPAS) is then the fraction of these transformed images that are classified in the adversary’s desired target class divided by the total number of transformed images.

**Perceptibility and confidence:** The perceptibility of the perturbations is measured by computing the average  $L_1$  norm of the perturbation for all adversarial examples. We compute and report the average confidence of the target class over all *successful* adversarial examples.

**Efficiency:** The slowest step in the generation of each adversarial example is running the optimization from Eq. (2). Each example takes about 60s to generate on the hardware we use. We stop the optimization run after 3000 steps for each sample since the change in loss between subsequent runs is vanishingly small.

**Remark (virtual and physical settings):** All results reported in the *virtual* setting involve benign and adversarial examples that remain in a digital format throughout. The results for *simulated physical attacks* (SPA) also pertain to images of this nature. The evaluation setup and metrics used for *real-world attacks* are described in Section 6.3.

## 6 WHITE-BOX ATTACKS

In this section we present results for both virtual and real-world attacks, in the white-box setting for adversarial examples. We also evaluate the effectiveness of our attacks against classifiers trained using adversarial training and its variants [14, 51].

**Attack pipeline in practice:** Equation 2 has a number of parameters such as batch-size  $B$ ,  $K$  to adjust the confidence of the generated adversarial examples,  $c$  to control the trade-off between distance and loss etc. We tuned these parameters as described in Appendix E to achieve high SPAS rates. Our attack pipeline, described in Figure 4 up-samples the generated perturbations to match the dimensions of the original image. When testing, we again down-sample these images so that they can be classified.

### 6.1 Virtual Out-of-Distribution attacks

We first evaluate the effectiveness of the Out-of-Distribution attack on high-resolution images. This allows us to pick the adversarial examples that work well in this setting for the real-world experiments. We experiment with 2 types of Out-of-Distribution images: Logo and Custom Sign.

**6.1.1 Logo attack.** In this attack, we modify high-resolution images of innocuous signs and logos such that they are classified in the desired target class. Each of the 7 original logos is passed through our attack pipeline to create 20 different Logo adversarial examples, each meant to be classified as a different, randomly chosen target traffic signs, giving a total of 140 adversarial examples. This attack achieves an impressive **VAS of 85.71%** and **SPAS of 65.07%** with an average confidence of 0.975 as reported in Table 2. While the adversarial perturbation on the logos significantly affects the classification results, it is generally not very visible to humans and usually blends in with the detail of the logos. Some successful Logo attacks are displayed in Figures 5a and 14b (Appendix).

**6.1.2 Custom Sign attack.** 3 circular blank signs of the colors blue, orange, and white are drawn on the computer and used as the original images for the Custom Sign attacks. Each blank sign is matched with six different masks (a cricle, a heart, a star, a cross, the phrase "Hello, World!" and a combination of these shapes). Each of the total 18 pairs of a blank sign and a mask is used to generate 10 Custom Sign attacks which would be classified as 10 randomly chosen target traffic signs. In total, 180 Custom Sign attacks are created for the virtual evaluation. Some Custom Sign attacks are shown in Figures 5b and 14c (Appendix). The attack produces adversarial signs that contain the shape of the mask filled with various colors. In Figure 5b, we pick some of the signs whose mask is filled well so the text or the shape is clearly visible. Some of the attacks do not fill up the entire mask resulting in incomplete shapes. This attack is considerably different from the Logo and the In-Distribution attacks as the optimization constraint is moved from the norm of the perturbation to its location instead. The adversary is allowed to add arbitrary perturbations as long as they are within the mask, hence the average  $L_1$  norm is not applicable. This new constraint seems to be more strict than the previous one, causing the attack success rate to drop as shown in Table 2.

**Main takeaway:** Both the Out-of-Distribution Logo and Custom Sign attacks are feasible in the virtual white-box setting with high-confidence classification for the successful adversarial examples. Their non-trivial SPAS enables the adversary to pick effective adversarial examples from this setting for real-world attacks.

### 6.2 Virtual In-Distribution attacks

We use high-resolution images from the auxiliary traffic data to generate In-Distribution attacks for the physical setting since the GTSRB test images are too low-resolution for large-scale printing. Both classifiers achieve 100% classification accuracy on these images in the benign setting. Similar to the Logo attack experiment, each of the 22 images in the auxiliary traffic data is used to generate 20 In-Distribution attacks for 20 randomly chosen target traffic signs. Therefore, one experiment contains 440 attacks in total. From

Attacks	White-box (Avg. confidence)	Black-box (Avg. confidence)
In-Distribution (auxiliary traffic data)	92.82% (0.9632)	96.68% (0.9256)
Out-of-Distribution (Logo)	52.50% (0.9524)	32.73% (0.9172)
Out-of-Distribution (Custom Sign)	96.51% (0.9476)	97.71% (0.9161)

**Table 3: Real-world attack success rates against the Multi-scale CNN in the white-box setting (average of 2 signs for each attack) and on Standard CNN in the black-box setting (best performing sign samples transferred from Multi-scale CNN).**

Table 2, we can see that attacks using these samples are feasible, confirming previous work.

### 6.3 Real-world attacks (Out-of-Distribution and In-Distribution)

To carry out a real-world attack, we first chose 2 adversarial examples from each of the Logo, Custom Sign and Auxiliary traffic sign datasets that performed well under the Simulated Physical Attack (i.e. had consistent high confidence classification). Each of these examples was re-sized to  $30 \times 30$  inches and printed on a high-quality poster. The printed signs are stuck on a pole at a height of 2 meters from the ground on the left side of the road with respect to a car driving towards the pole. A GoPro HERO5 was mounted behind the car’s windshield to take videos of  $2704 \times 1520$  pixels at 30 frames per second. Starting around 25 meters from the sign, the car approached it with an approximate speed of 16kph. The traffic sign recognition pipeline only detects and classifies once every 5 frames to reduce the processing time and redundancy.

**Real-world attack metric:** To evaluate real-world attacks, we use a different metric since the number of adversarial examples used is smaller. In this section, the attack success rate reported is computed by counting the number of frames in which a sign was detected in the input frame and classified as the target class and dividing this number by the total number of frames in which there was a detection. In other words, we can express this as:

$$\text{Drive-by attack success} = \frac{\text{No. of frames sign is misclassified}}{\text{No. of frames sign is detected}} \quad (4)$$

**Real-world attack success:** To demonstrate the effectiveness of our adversarial examples in the real world, we carried out *drive-by tests* (shown in Figure 8) on 2 samples from each of our attacks (In-Distribution, Logo, and Custom Sign). Each of the drive-by attack success rates reported in Table 3 is an average of three runs. We hypothesize that the lower attack success rate for the Logo attack in the drive-by setting occurs because the particular instance that was chosen had a more difficult target class. However, for frames classified as the target, the confidence is high. The source-target pairing for the In-Distribution attack was two speed limit signs, which could have contributed to the high attack success.

**Main takeaway.** In the drive-by tests, both the Out-of-Distribution and In-Distribution adversarial examples can achieve attack success rates in excess of 90%.

Attacks	VAS	SPAS	Avg. confidence	Avg. norm ( $L_1$ )
In-Distribution (auxiliary traffic data)	2.53%	2.47%	0.9358	38.92
Out-of-Distribution (Logo)	11.42%	7.42%	0.8054	36.12
Out-of-Distribution (Custom Sign)	6.67%	5.77%	0.9957	N.A.

**Table 4: Attack success rates and deterioration rate on adversarially trained Multi-scale CNN for Out-of-Distribution and In-Distribution attacks in the virtual white-box setting.**

### 6.4 Attacking defended models

In this section we examine the effectiveness of the Out-of-Distribution and In-Distribution attacks against state-of-the-art defenses based on the concept of *adversarial training* [14]. In this paper, we are the first to analyze the effectiveness of adversarial training against physically robust adversarial examples. We note that this defense mechanism is completely ineffective against lenticular printing based attacks.

**6.4.1 Adversarial training.** This defense modifies the loss function of neural networks to increase their robustness against adversarial examples. The *training loss* is modified as follows:

$$\ell_f^{\text{adv}}(\mathbf{x}, \mathbf{y}) = \alpha \ell_f(\mathbf{x}, \mathbf{y}) + (1 - \alpha) \ell_f(\tilde{\mathbf{x}}, \mathbf{y}), \quad (5)$$

where  $\alpha \in [0, 1]$  controls the adversarial component of the loss and  $\tilde{\mathbf{x}}$  is an adversarial example. We trained Multi-scale CNN using the loss function above with  $\alpha = 0.5$  and using Fast Gradient Sign adversarial examples [14] with an  $\epsilon = 0.3$ . The model was trained for 15 epochs and has an accuracy of **96.37%** on the GTSRB validation set. We do not use the adversarial examples generated using Equation (2) as each sample takes prohibitively long to generate for the purposes of training. Further details are in Appendix G.

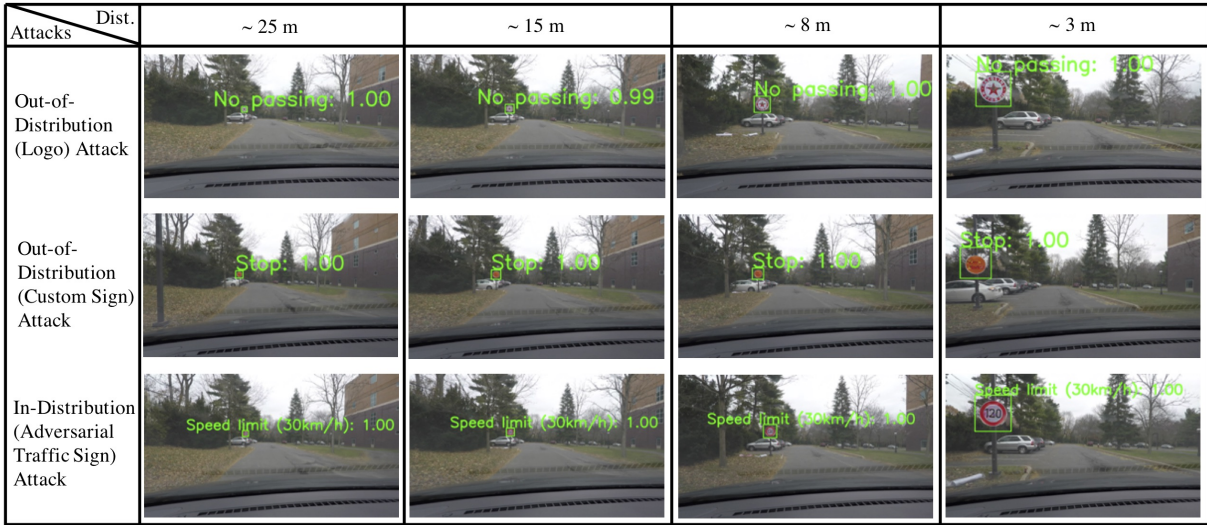
**6.4.2 Effect of Out-of-Distribution attacks on adversarially trained models.** The virtual white-box attack results given in Table 4 demonstrate that adversarial training is very effective as a defense against In-Distribution attacks since the VAS using auxiliary traffic data drops to 54.34% to 2.53% while the SPAS drops from 36.65% to 2.47%. On the other hand, the VAS and SPAS for the Out-of-Distribution Logo attack remain at 11.42% and 7.42% respectively. Thus, the Logo attack is about 3 – 4× more effective than the In-Distribution attack. The Custom Sign attack is also more effective than the In-Distribution attack, indicating that Out-of-Distribution attacks are more effective against these defenses.

**Main takeaway:** Adversarial training cannot defend against Lenticular Printing and Out-of-Distribution attacks are more effective than previously examined In-Distribution attacks. Thus, we have introduced new attack vectors against these state-of-the-art defenses.

## 7 BLACK-BOX ATTACKS ARE POSSIBLE

In the black-box attack setting, we use adversarial examples generated for the Multi-scale CNN and test their attack success rate on Standard CNN in both virtual and real-world settings. The recognition pipeline is kept the same, but the classifier is changed. The reason this attack is expected to work is the well-known phenomenon





**Figure 8: Classification outcomes for adversarial examples in the drive-by test. All combinations of distances and attacks lead to the desired classification outcome with high confidence. Even frames captured at around 25 metres from the sign lead to high confidence targeted misclassification.**

Attacks	Black-box VAS	Black-box SPAS	Avg. confidence
In-Distribution (auxiliary traffic data)	7.14%	6.08%	0.9273
Out-of-Distribution (Logo)	14.28%	12.57%	0.8717
Out-of-Distribution (Custom Sign)	3.33%	6.00%	0.7193

**Table 5: Black-box attack success rates for In-Distribution and Out-of-Distribution attacks. Adversarial examples generated for Multi-scale CNN are tested on Standard CNN.**

of *transferability* [18–20], where adversarial examples generated for one classifier remain adversarial for another. Since both classifiers are trained on the same dataset, we are implicitly assuming that the adversary has access to the training data. This is a common assumption in previous work on transferability-based black-box attacks.

**Virtual attacks:** Table 5 shows adversarial success rate in a black-box setting, demonstrating the transferability of adversarial examples. While there is a significant drop in attack success rates, both the VAS and SPAS for the Out-of-Distribution Logo attacks remain at non-negligible rates of above 10% indicating that even an attacker with very limited knowledge of the target classifier can carry out successful attacks.

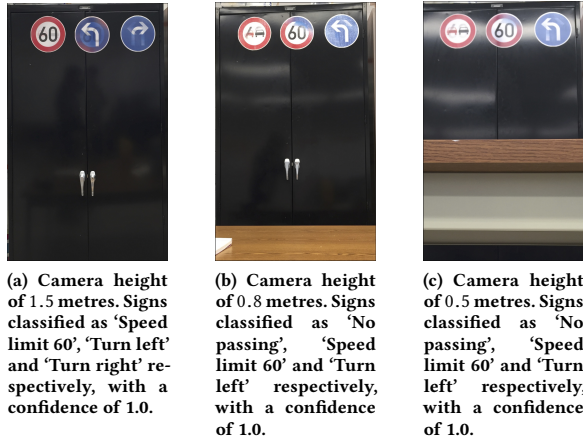
**Real-world attacks:** We use the black-box attack pipeline to evaluate the same videos captured in the drive-by tests (Table 3). With the same set of signs that perform best in the white-box case, the Out-of-Distribution Custom Sign and the In-Distribution attacks both achieve very attack high success rates of **97.71%** and **96.68%** respectively, comparable to the white-box setting. One of the Logo attacks achieves an attack success rate of 32.73% in the black-box case. Note that the higher attack success rates for the black-box setting arise from the fact that we just report numbers for the best-performing adversarial sign and not an average over two signs. **Main takeaway:** Black-box attack success rates in excess of 90% can be achieved in real-world settings with appropriately chosen

adversarial examples. This represents a significant threat against traffic sign recognition systems since no system knowledge is required for these.

## 7.1 Lenticular printing attacks

In these experiments, we aim to show that the difference in viewing angle between the camera on an AV and the human controller or passenger can lead to an attack, since the human will see an appropriate sign, while the camera will see the targeted sign, which could be dangerous in that environment. To simulate what the human/camera will see, we take pictures from different heights showing the sign changing. We emphasize that the adversarial nature of these examples stems purely from the viewing angle, and the human and the camera would recognize the sign as the same if they viewed it from the same angle.

In our experiments with the Lenticular Printing attack, we stuck three different signs on an indoor surface (Figure 9). These signs flip between ‘Speed limit 60’ and ‘No Passing’, ‘Turn left’ and ‘Speed limit 60’, and ‘Turn left’ and ‘Turn right’ respectively. We take pictures of these signs from three different observation angles ( $\theta_1 = 8$ ,  $\theta_2 = 25$ , and  $\theta_3 = 33$  degrees) using a standard mobile phone camera at a distance of 1.9 metres from the signs. The camera is held parallel to the plane of the cupboard. In all cases, when passed through our traffic sign recognition pipeline, each sign classified was classified as the appropriate one for that viewing angle with high confidence. **Validating theory:** Based on the theoretical discussion in Section 4, we expect that the camera observes two images from two different angles if the difference between these two observation angles is between 12.5 and 25 degrees. This matches our observations demonstrated in Figure 9: (1)  $25 \geq |\theta_1 - \theta_2| = 17 \geq 12.5$  and  $25 \geq |\theta_1 - \theta_3| = 25 \geq 12.5$ , indicating that the first image should look different from other images, and (2)  $|\theta_2 - \theta_3| = 8 \leq 12.5$ , suggesting that the second and third images should look similar.



**Figure 9: Proof-of-concept implementation of the Lenticular Printing attack.** These images show that if the camera used for the traffic sign recognition module of an AV is at a different height from the human controller, then the Lenticular Printing attack can fool the classifier, while appearing to be correct to the human.

**How can the attacker ensure that the driver cannot see different images as the car approaches the sign?** The attacker can readily measure the range of the observation angle of the driver and the camera given the height of the driver’s eyes, the height of the camera, and height of the sign. The attacker can then create and install the sign, while considering the theoretical discussion in Section 4. A straightforward estimation of the height of the driver’s eyes and the height of the camera can be made based on the height of the target vehicle. For simplicity, the attacker can install the sign on a flat road at the height of the driver’s eyes. In this case the observation angle of the driver remains fixed, in particular,  $\theta_1 = 0$ .

## 8 LIMITATIONS AND FUTURE WORK

In this section, we discuss some limitations of our approaches for both attacks and defenses and outline possibilities for future work to overcome these limitations.

**Adversarial example detectors as a countermeasure:** While detection based defenses such as Feature Squeezing [52] and Magnet [53] are ineffective against In-Distribution white-box attacks [54], it is an open question as to how they will perform against Out-of-Distribution attacks. We plan to explore detection-based defenses such as these in future work.

**Fooling synthesis of sensor inputs:** The computer vision subsystem of an AV, while critical, is not the only actor in the decision making process when confronted with a new situation. A number of other sensors such as LIDAR, GPS, radar etc. also provide inputs which are then synthesized to come up with a decision [55]. While the computer vision subsystem is the only one able to recognize a traffic sign, the other sensors may be able to indicate that the sign recognized is incompatible with their inputs. This consideration is out of the scope of the current work, but we plan to explore simultaneous attacks on these varied subsystems in future work.

**Choice of norm:** In this paper, we chose the  $L_1$  norm to measure the visibility of adversarial perturbations as done in several previous works on virtual attacks [15, 56, 57]. However, it is still an open research question if this is the best choice of distance function to constrain adversarial perturbations. Previous work on generating adversarial examples has explored the appropriateness of other commonly used Euclidean norms such as  $L_\infty$  [14] and  $L_2$  [15, 19] as well. Especially in the context of future work on generating physically realizable adversarial examples, we encourage further work on conducting user studies to determine the most appropriate proxy for measuring human perceptibility of adversarial perturbations.

## 9 RELATED WORK

**Virtual white-box attacks:** Evasion attacks (test phase) have been proposed for Support Vector Machines [11, 12], random forests [58] and neural networks [12–17]. Nguyen et al. [59] generate adversarial examples starting from random noise restricted to the virtual white-box setting. Poisoning attacks (training phase) have also been proposed for a variety of classifiers and generative models [60–64]. Attacks have also been proposed on policies for reinforcement learning [65], models with structured prediction outputs [66], semantic segmentation models [67, 68] and neural network based detectors [69].

**Virtual black-box attacks:** The phenomenon of transferability [13, 19] has been used to carry out black-box attacks on real-world models [18, 20]. Powerful query-based black-box attacks have also been demonstrated [21–23, 70–72].

**Real-world attacks:** Kurakin et al. [24] were the first to explore real-world adversarial examples. They printed out adversarial examples generated using the Fast Gradient Sign attack and passed them through a camera to determine if the resulting image was still adversarial. This attack was restricted to the white-box setting and the effect of varying physical conditions was not taken into account. Sharif et al. [44] investigated how face recognition systems could be fooled by having a subject wear eyeglasses with adversarial perturbations printed on them. While their attack results were encouraging, they did not rigorously account for varying physical conditions and only took multiple pictures of a subject’s face to increase the robustness of the attack. Further, they only tested their black-box attacks in a virtual setting with query access to the target model. Petit et al. [73] examine the susceptibility of the LIDAR and camera sensors in an autonomous car but their attacks are unable to cause targeted misclassification. Lu et al. [74] and Chen et al. [75] attempted physical-world attacks on R-CNN based traffic sign detectors [76]. However, large perturbations were needed in both cases. We plan to explore Out-of-Distribution attacks on these detectors in future work. Comparisons with the related work of Athlaye et al. [26] and Evtimov et al. [27] have been elucidated in Section 1.

## 10 CONCLUSION

In this paper, we have demonstrated a wide range of attacks on traffic sign recognition systems, which have severe consequences for self-driving cars. Out-of-Distribution attacks allow an adversary to convert any sign or logo into into a targeted adversarial example.

The Lenticular Printing attack moves beyond the paradigm of adversarial examples to create images that look different from varying heights, allowing an adversary to stealthily embed a potentially dangerous traffic sign into an innocuous one, with no access to the internals of the classifier. We demonstrated the effectiveness of our attacks in both virtual and real-world settings. We are the first to carry out black-box attacks in a real-world setting as well as to evaluate possible countermeasures against physical realizations of robust adversarial examples. We hope our discussion of future research directions encourages further exploration into securing physically deployed machine learning systems.

## REFERENCES

- [1] A. Mosenia, S. Sur-Kolay, A. Raghunathan, and N. K. Jha. Caba: Continuous authentication based on bioaura. *IEEE Transactions on Computers*, 66(5):759–772, May 2017.
- [2] SYMM Kung, M Mak, and S Lin. *Biometric authentication: a machine learning approach*. Prentice Hall Press, 2004.
- [3] A. Mosenia, S. Sur-Kolay, A. Raghunathan, and N. K. Jha. Wearable medical sensor-based system design: A survey. *IEEE Transactions on Multi-Scale Computing Systems*, 3(2):124–138, 2017.
- [4] A. M. Nia, M. Mozaffari-Kermani, S. Sur-Kolay, A. Raghunathan, and N. K. Jha. Energy-efficient long-term continuous personal health monitoring. *IEEE Transactions on Multi-Scale Computing Systems*, 1(2):85–98, 2015.
- [5] NVIDIA. Self driving vehicles development platform. <http://www.nvidia.com/object/drive-px.html>. Accessed: 2016-10-31.
- [6] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- [7] A. Mosenia, J. F. Bechara, T. Zhang, P. Mittal, and M. Chiang. Promotive: Bringing programability and connectivity into isolated vehicles. *Accepted for publication in Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies (IMWUT)*, 2018.
- [8] A. Mosenia, X. Dai, P. Mittal, and N. Jha. Pinme: Tracking a smartphone user around the world. *IEEE Transactions on Multi-Scale Computing Systems*, PP(99):1–1, 2017.
- [9] M. Mozaffari-Kermani, S. Sur-Kolay, A. Raghunathan, and N. K. Jha. Systematic poisoning attacks on and defenses for machine learning in healthcare. *IEEE Journal of Biomedical and Health Informatics*, 19(6):1893–1905, Nov 2015.
- [10] Arsalan Mosenia and Niraj Jha. A comprehensive study of security of Internet of Things. *IEEE Trans. Emerging Topics in Computing*, 5(4):586–602, 2017.
- [11] Battista Biggio, Igino Corona, Blaine Nelson, Benjamin IP Rubinstein, Davide Maiorca, Giorgio Fumera, Giorgio Giacinto, and Fabio Roli. Security evaluation of support vector machines in adversarial environments. In *Support Vector Machines Applications*, pages 105–153. Springer, 2014.
- [12] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. Deepfool: a simple and accurate method to fool deep neural networks. *arXiv preprint arXiv:1511.04599*, 2015.
- [13] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- [14] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- [15] Nicholas Carlini and David A. Wagner. Towards evaluating the robustness of neural networks. *CoRR*, abs/1608.04644, 2016.
- [16] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. *arXiv preprint arXiv:1610.08401*, 2016.
- [17] Nicolas Papernot, Patrick McDaniel, Somesh Jha, Matt Fredrikson, Z Berkay Celik, and Ananthram Swami. The limitations of deep learning in adversarial settings. In *2016 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 372–387. IEEE, 2016.
- [18] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Ananthram Swami. Practical black-box attacks against deep learning systems using adversarial examples. In *Proceedings of the 2017 ACM Asia Conference on Computer and Communications Security*, 2017.
- [19] Nicolas Papernot, Patrick McDaniel, and Ian Goodfellow. Transferability in machine learning: from phenomena to black-box attacks using adversarial samples. *arXiv preprint arXiv:1605.07277*, 2016.
- [20] Yanpei Liu, Xinyun Chen, Chang Liu, and Dawn Song. Delving into transferable adversarial examples and black-box attacks. In *ICLR*, 2017.
- [21] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*, 2017.
- [22] Arjun Nitin Bhagoji, Warren He, Bo Li, and Dawn Song. Exploring the space of black-box attacks on deep neural networks. *arXiv preprint arXiv:1712.09491*, 2017.
- [23] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 15–26. ACM, 2017.
- [24] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*, 2016.
- [25] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, pages 1528–1540. ACM, 2016.
- [26] Anish Athalye, Logan Engstrom, Andrew Ilyas, and Kevin Kwok. Synthesizing robust adversarial examples. *CoRR*, abs/1707.07397, 2017.
- [27] Ivan Evtimov, Kevin Eykholt, Earlene Fernandes, Tadayoshi Kohno, Bo Li, Atul Prakash, Amir Rahmati, and Dawn Song. Robust physical-world attacks on machine learning models. *CoRR*, abs/1707.08945, 2017.
- [28] Tesla. Autopilot | tesla. <https://www.tesla.com/autopilot>. Accessed: 2017-12-05.
- [29] James Vincent. Apple’s latest ai research explores the problem of mapping systems for self-driving cars. <https://www.theverge.com/2017/11/22/16689810/apple-ai-research-self-driving-cars-autonomous>. Accessed: 2017-12-10.
- [30] Chawin Sitawarin, Arjun Nitin Bhagoji, Arsalan Mosenia, Mung Chiang, and Prateek Mittal. Rogue signs: Deceiving traffic sign recognition with malicious ads and logos. In *1st Deep Learning and Security Workshop (IEEE S&P 2018)*, DLS 2018, 2018.
- [31] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [32] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [33] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [34] Yaniv Taigman, Ming Yang, Marc’Aurelio Ranzato, and Lior Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1701–1708, 2014.
- [35] Vincenzo Barrile, Giuseppe M Meduri, and Domenico Cuzzocrea. Automatic recognition of road signs by hough transform: Road-gis. *Journal of Earth Science and Engineering*, 2(1), 2012.
- [36] P. Yakimov and V. Fursov. Traffic signs detection and tracking using modified hough transform. In *2015 12th International Joint Conference on e-Business and Telecommunications (ICETE)*, volume 05, pages 22–28, July 2015.
- [37] Miguel Ángel García-Garrido, Miguel Ángel Sotelo, and Ernesto Martín-Gorostiza. *Fast Road Sign Detection Using Hough Transform for Assisted Driving of Road Vehicles*, pages 543–548. Springer Berlin Heidelberg, Berlin, Heidelberg, 2005.
- [38] Hough circle transform – opencv-python tutorials 1 documentation. [http://opencv-python-tutorials.readthedocs.io/en/latest/py\\_tutorials/py\\_imgproc/py\\_houghcircles/py\\_houghcircles.html](http://opencv-python-tutorials.readthedocs.io/en/latest/py_tutorials/py_imgproc/py_houghcircles/py_houghcircles.html). (Accessed on 12/11/2017).
- [39] J. Canny. A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-8(6):679–698, Nov 1986.
- [40] Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 387–402. Springer, 2013.
- [41] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- [42] Florian Tramèr, Fan Zhang, Ari Juels, Michael K Reiter, and Thomas Ristenpart. Stealing machine learning models via prediction apis. In *USENIX Security*, 2016.
- [43] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [44] Mahmood Sharif, Sruti Bhagavatula, Lujo Bauer, and Michael K. Reiter. Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, CCS ’16, pages 1528–1540, New York, NY, USA, 2016. ACM.
- [45] Image segmentation – skimage v0.14dev docs. [http://scikit-image.org/docs/dev/user\\_guide/tutorial\\_segmentation.html](http://scikit-image.org/docs/dev/user_guide/tutorial_segmentation.html). (Accessed on 12/12/2017).
- [46] John Canny. A computational approach to edge detection. In *Readings in Computer Vision*, pages 184–203. Elsevier, 1987.
- [47] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 2012.
- [48] Sebastian Houben, Johannes Stallkamp, Jan Salmen, Marc Schlipsing, and Christian Igel. Detection of traffic signs in real-world images: The German Traffic



- Sign Detection Benchmark. In *International Joint Conference on Neural Networks*, 2013.
- [49] P. Sermanet and Y. LeCun. Traffic sign recognition with multi-scale convolutional networks. In *The 2011 International Joint Conference on Neural Networks*, pages 2809–2813, July 2011.
- [50] freesuperflip. <http://www.vuethru.com/freesuperflip.html>. (Accessed on 02/13/2018).
- [51] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv:1706.06083 [cs, stat]*, June 2017.
- [52] Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. *arXiv preprint arXiv:1704.01155*, 2017.
- [53] Dongyu Meng and Hao Chen. Magnet: a two-pronged defense against adversarial examples. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 135–147. ACM, 2017.
- [54] Nicholas Carlini and David Wagner. Magnet and<sup>2</sup> efficient defenses against adversarial attacks<sup>3</sup> are not robust to adversarial examples. *arXiv preprint arXiv:1711.08478*, 2017.
- [55] Fernando Mujica. Scalable electronics driving autonomous vehicle technologies. *Texas Instrument*, 2014.
- [56] Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Analysis of classifiers’ robustness to adversarial perturbations. *arXiv preprint arXiv:1502.02590*, 2015.
- [57] Pin-Yu Chen, Yash Sharma, Huan Zhang, Jinfeng Yi, and Cho-Jui Hsieh. Ead: elastic-net attacks to deep neural networks via adversarial examples. *arXiv preprint arXiv:1709.04114*, 2017.
- [58] Alex Kantchelian, JD Tygar, and Anthony D Joseph. Evasion and hardening of tree ensemble classifiers. In *Proceedings of the 33rd International Conference on Machine Learning (ICML-16)*, 2016.
- [59] Anh Nguyen, Jason Yosinski, and Jeff Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 427–436. IEEE, 2015.
- [60] Battista Biggio, Blaine Nelson, and Pavel Laskov. Poisoning attacks against support vector machines. In *Proceedings of the 29th International Conference on Machine Learning (ICML-12)*, pages 1807–1814, 2012.
- [61] Benjamin IP Rubinstein, Blaine Nelson, Ling Huang, Anthony D Joseph, Shing-hon Lau, Satish Rao, Nina Taft, and JD Tygar. Stealthy poisoning attacks on pca-based anomaly detectors. *ACM SIGMETRICS Performance Evaluation Review*, 37(2):73–74, 2009.
- [62] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is feature selection secure against training data poisoning? In *ICML*, 2015.
- [63] Shike Mei and Xiaojin Zhu. The security of latent dirichlet allocation. In *AISTATS*, 2015.
- [64] Scott Alfeld, Xiaojin Zhu, and Paul Barford. Data poisoning attacks against autoregressive models. In *AAAI*, 2016.
- [65] Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *arXiv preprint arXiv:1702.02284*, 2017.
- [66] Moustapha Cisse, Yossi Adi, Natalia Neverova, and Joseph Keshet. Houdini: Fooling deep structured prediction models. *arXiv preprint arXiv:1707.05373*, 2017.
- [67] Anurag Arnab, Ondrej Miksik, and Philip HS Torr. On the robustness of semantic segmentation models to adversarial attacks. *arXiv preprint arXiv:1711.09856*, 2017.
- [68] Volker Fischer, Mummadi Chaithanya Kumar, Jan Hendrik Metzen, and Thomas Brox. Adversarial examples for semantic image segmentation. *arXiv preprint arXiv:1703.01101*, 2017.
- [69] Cihang Xie, Jianyu Wang, Zhishuai Zhang, Yuyin Zhou, Lingxi Xie, and Alan Yuille. Adversarial examples for semantic segmentation and object detection. In *International Conference on Computer Vision*. IEEE, 2017.
- [70] Blaine Nelson, Benjamin IP Rubinstein, Ling Huang, Anthony D Joseph, Shing-hon Lau, Steven J Lee, Satish Rao, Anthony Tran, and JD Tygar. Near-optimal evasion of convex-inducing classifiers. In *AISTATS*, pages 549–556, 2010.
- [71] Blaine Nelson, Benjamin IP Rubinstein, Ling Huang, Anthony D Joseph, Steven J Lee, Satish Rao, and JD Tygar. Query strategies for evading convex-inducing classifiers. *The Journal of Machine Learning Research*, 13(1):1293–1332, 2012.
- [72] Daniel Lowd and Christopher Meek. Adversarial learning. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 641–647. ACM, 2005.
- [73] Jonathan Petit, Bas Stottelaar, Michael Feiri, and Frank Kargl. Remote attacks on automated vehicles sensors: Experiments on camera and lidar. *Black Hat Europe*, 11:2015, 2015.
- [74] Jiajun Lu, Hussein Sibai, and Evan Fabry. Adversarial examples that fool detectors. *arXiv preprint arXiv:1712.02494*, 2017.
- [75] Shang-Tse Chen, Cory Cornelius, Jason Martin, and Duen Horng Chau. Robust physical adversarial attack on faster r-cnn object detector. *arXiv preprint arXiv:1804.05810*, 2018.
- [76] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *CoRR*, abs/1506.01497, 2015.
- [77] A. Mogelmoose, M. M. Trivedi, and T. B. Moeslund. Vision-based traffic sign detection and analysis for intelligent driver assistance systems: Perspectives and survey. *IEEE Transactions on Intelligent Transportation Systems*, 13(4):1484–1497, Dec 2012.
- [78] Traffic signs classification with a convolutional network - alex staravoitau’s blog. <https://navoshta.com/traffic-signs-classification/>. (Accessed on 12/12/2017).
- [79] François Chollet et al. Keras. <https://github.com/keras-team/keras>, 2015.
- [80] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [81] Mark Everingham, Luc Van Gool, Christopher K. I. Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, Jun 2010.

## A APPENDIX: DATASET DETAILS

We use the German Traffic Sign Recognition Benchmark (GTSRB) [47], a widely-used standard dataset for traffic sign recognition, to train and test both classifiers<sup>3</sup>. It contains 51,839 images of 43 types of traffic signs in total. Each image is a  $32 \times 32$  RGB image with the pixels scaled to lie in  $[0, 1]$ . Sample images are shown in Figure 10. The 39,209 training samples are augmented to 86,000 (2,000 for each class) using perspective transformation and flipping signs that are invariant to flipping (i.e. "no entry") or that are interpreted as a different class after flipped (i.e. "turn right" to "turn left").

**Auxiliary high-resolution dataset:** To create fake traffic signs that look realistic, we used 22 high-resolution traffic sign images to use as original images for the In-Distribution attack, shown in Figure 11. They are a mixture of real photographs and computer-generated drawings on an arbitrary background image.

**Logo dataset:** To create realistic-looking Out-of-Distribution Logo adversarial examples, we used 7 high-resolution Logo images on arbitrary backgrounds as shown in Figure 12.

**Custom Sign dataset:** To create realistic-looking Out-of-Distribution Custom Sign adversarial examples, we used 3 circular blank signs (orange, blue and white) in combination with 6 masks on arbitrary backgrounds as shown in Figure 13.

## B APPENDIX: DETECTOR DETAILS

For simplicity and without loss of generality, we design our detector to only locate circular signs on images using a well-known shape-based object detection technique in computer vision, Hough transform [38]. Triangular signs can be detected by a similar method described in [36]. In detail, our pipeline takes a video frame as an input and smooths it with a Gaussian filter to remove random noise. The processed images are passed through a Canny edge detector and then a circle Hough transform which outputs coordinates of the center and radii of the circles detected. The outputs are used to determine a square bounding box around the detected signs. The section of the input frame inside the bounding boxes are cropped out of their original *unprocessed* image and resized to  $32 \times 32$  pixels which is the input size of the neural network classifier.

<sup>3</sup>We chose GTSRB over the LISA Traffic Sign Dataset [77] as GTSRB offers a much larger number of samples leading to better generalized classifiers.



Figure 10: Samples from the GTSRB dataset



Figure 11: The auxiliary traffic data. Some high-resolution images sampled from the auxiliary traffic data. The auxiliary traffic data is a combination of real traffic sign photos and computer-generated drawings in front of a simple background.



Figure 12: Logo dataset. High-resolution images of common logos.

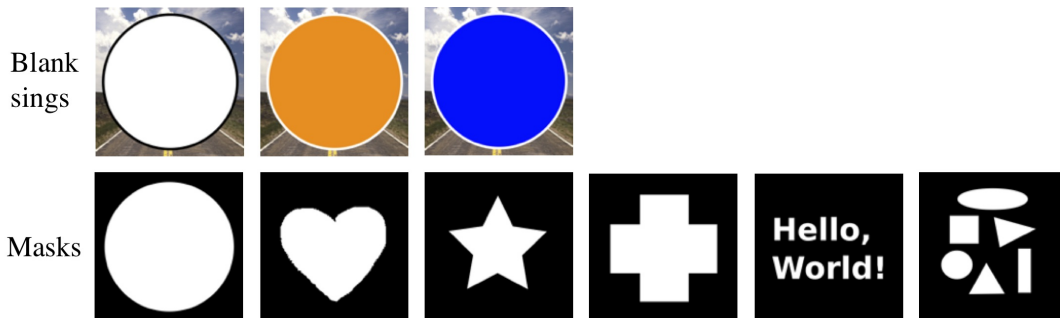


Figure 13: Custom Sign dataset. Blank signs with 6 different masks.

### C APPENDIX: CLASSIFIER DETAILS

Our classifier in the *white-box setting*, referred to as Multi-scale CNN is based on a multi-scale convolutional neural network (CNN) [49] with 3 convolutional layers and 2 fully connected layers. There is also an extra layer that pools and concatenates features from all convolutional layers, hence the name multi-scale. The training procedure is adapted from [78]. The model is trained on data-augmented

training set (86,000 samples, 2,000 for each class) generated by random perspective transformation and random brightness and color adjustment. There is no image preprocessing. The training takes 41 epochs with learning rate of 0.001. Dropout and weight regularization ( $\lambda = 1e-3$ ) are used to reduce overfitting. The accuracy of the model, which will be referred to as Multi-scale CNN, on the validation set is **98.50%**. The model is defined in Keras [79] using Tensorflow [80] backend. The entire recognition pipeline (the

detector combined with the classifier) is tested with the German Traffic Sign Detection Benchmark (GTSDB) [48] where it achieves a reasonable mean average precision (mAP) [81] of 77.10% at the intersection-over-union (IoU) of 0.5.

Standard CNN (used in the black-box attack setting) is a standard CNN (not multi-scale) with four convolutional layers and three dense layers. It is trained on the same dataset as Multi-scale CNN, and achieves an accuracy of 98.66% on the validation set of GTSRB. The mean average precision (mAP) of the entire pipeline with Standard CNN on the GTSDB dataset is 81.54%. In the black-box setting, the attacks are generated based on Multi-scale CNN which is assumed to be trained by the adversary as a substitute model for the true classifier in use, Standard CNN. All experiments in Section 7 use this black-box pipeline to evaluate the attacks. Both classifiers achieve a 100% accuracy on the auxiliary traffic data, implying that they have generalized well to images of traffic signs.

## D APPENDIX: BASELINE RESULTS

We compare the adversarial examples generated using our method to those generated by the optimization problem in Eq. (1) (referred to as Vanilla Optimization) on a random subset of 1000 traffic signs chosen from the test data of the GTSRB dataset. These results are not comparable to the ones in the main text using high-resolution data since both the masking and re-sizing effects are absent from these.

The result shown in Table 6 demonstrates that our attack has a much lower deterioration rate (DR) compared to the Vanilla Optimization attacks. The Vanilla Optimization attack with the same parameters as our attack indicates that forcing perturbation to be large (having a large norm) can increase its robustness to transformation to a certain degree. However, our attack can achieve a much lower deterioration rate with a similar average norm emphasizing the substantial effect the addition of random transformations to the optimization problem has on improving the physical robustness of adversarial examples. As expected, both the VAS and SPAS are lower for the adversarially trained Multi-scale CNN and the black-box attack setting.

## E APPENDIX: PARAMETER TUNING

We use an additional metric, the *Deterioration rate (DR)* in this section. This metric measures what fraction of the generated adversarial examples remain adversarial under simulated physical conditions:  $DR = 1 - \frac{SPAS}{VAS}$ . A higher deterioration rate implies that the adversarial examples are more likely to degrade under the transformations, and by extension in a real-world setting.

The set of initial parameters was manually chosen to maximize the VAS on the auxiliary traffic data on Multi-scale CNN while maintaining a low average perturbation as measured in the  $L_1$  norm. We use the  $L_1$  norm since we found it to achieve the best trade-off between the various performance metrics such as VAS, SPAS and DR that we considered. Due to this trade-off, we chose not to use a grid search to fine the optimization parameters. Once we had a set of parameters that worked well, we changed each of the parameters individually to understand their effect on the attack performance metrics. First, we set the parameter  $L$  to control the norm of the output perturbation. Depending on the type of

norms used in the objective function,  $L$  can be chosen roughly to any number larger than zero to force the algorithm to explore more solutions that cause targeted misclassification instead of finding those that merely minimize the norm. With  $L$  chosen for each norm, we vary the constant  $c$ , norm  $p$ , the number of transformation  $T$ , and degree of randomness of the transformations. We use  $\theta_1$  for perspective transformation,  $\theta_2$  for brightness adjustment and  $\theta_3$  for image resizing. We find that varying  $\theta_3$  does not significantly affect the outputs and thus do not discuss it further.

The baseline attack uses  $L_1$  norm with the parameters  $c = 3$ ,  $K = 100$ ,  $L = 30$ ,  $T = 128$ ,  $\theta_1 = 0.07$ ,  $\theta_2 = 0.15$ . Table 7 shows results from 13 experiments each of which vary specified parameters from the baseline. We use Adam optimizer to solve the optimization problem with learning rate (step size) of 0.02 without decay. It must be noted that all four result columns of Table 7 must be viewed in conjunction, instead of only considering one column, because they represent the trade-offs of the attacks. For example, we can see that by increasing the degree of randomness ( $\theta_1, \theta_2$ ), attack success rates generally increase, but the norms also become larger making the perturbation more noticeable. In particular, setting  $\theta_1 = 0.15$  results in the second highest physical success rate and the lowest deterioration rate, but it also produces perturbation with the second largest norm. Similarly, setting  $K$  to be larger encourages the optimization to look for more *adversarial* solutions which also comes with the cost of a large perturbation. Increasing the number of transformations makes the attack more successful under transformation while reducing the success rate in the virtual setting slightly. In fact, using 512 transformations produces both high physical attack success rate and a small norm for the perturbation with only the expense of longer computation time. For an adversary who only needs a few adversarial signs, he or she can afford to use a large number of transformation or even run a search for optimal parameters for a specific setting. For the Custom Sign attack, since there is no constraint on the norm of the perturbation both  $c$  and  $L$  are increased to obtain perturbations that fit within the mask and are adversarial. Table 7 contains the results from tuning the various optimization parameters from Equation 2.

**Main takeaway.** The number of tunable parameters in the optimization represents an advantage for the adversary since they can be tuned to achieve the desired trade-off between the different performance metrics. We chose a particular set of parameters that worked well for our evaluation setting, and clarified the changes in performance that occur when these are tweaked.

## F APPENDIX: MORE ADVERSARIAL EXAMPLES

Here, we include some of the adversarial examples we have generated for all of the three attacks in Figure 14. The samples shown also achieve 100% SPAS.

## G APPENDIX: ADVERSARIAL TRAINING

Ideally, the  $\tilde{x}$  in Eq.(5) (modified loss function) should also be generated in the same manner as the robust physical attacks. However, as mentioned earlier, each adversarial example generated using our method takes around 60s to generate. This makes it impractical to run the optimization for every training batch.

Models	Attacks	Virtual attack success (VAS)	Simulated physical attack success (SPAS)	Avg. norm ( $L_1$ )	Avg. confidence
Multi-scale CNN	Baseline CW [15] attack (GTSRB test data)	97.91%	46.74%	30.45	0.9321
	In-Distribution (GTSRB test data)	99.07%	95.50%	31.43	0.9432
Multi-scale CNN <sub>adv</sub>	In-Distribution (GTSRB test data)	36.35%	27.52%	31.57	0.9428
Multi-scale CNN → Standard CNN	In-Distribution (GTSRB test data)	47.77%	38.08%	31.43	0.8838

**Table 6: White-box attack success rates for baseline and In-Distribution attacks on the GTSRB test data. The parameters for the Carlini-Wagner attack are modified in order to increase its Simulated Physical Attack Success rate.**

Description	Parameters	VAS	SPAS	DR	Avg. norm ( $L_1$ )
Chosen params	$c = 3, K = 100, L = 30, T = 128, \theta_1 = 0.07, \theta_2 = 0.15$	54.34%	36.65%	32.55%	37.71
Perturbation norm	$L_2$ norm ( $c = 0.2, L = 2$ )	27.04%	15.09%	37.23%	76.74
	$L_2$ norm ( $c = 0.02, L = 2$ )	14.03%	7.76%	44.69%	54.13
	$L_\infty$ norm ( $c = 5e-5, L = 0.1$ )	43.37%	24.41%	43.72%	162.76
Adversarial confidence	$K = 50$	48.21%	28.58%	40.71%	32.89
	$K = 200$	<b>59.69%</b>	<b>50.91%</b>	<b>14.71%</b>	<b>50.53</b>
# of transformed samples	$T = 32$	54.59%	32.08%	41.23%	35.44
	$T = 512$	46.43%	37.81%	18.57%	35.82
Degree of transformation ( $\theta_1$ : perspective transformation, $\theta_2$ : brightness adjustment)	$\theta_1 = 0, \theta_2 = 0$	31.89%	6.26%	80.37%	31.87
	$\theta_1 = 0.03$	44.13%	8.88%	79.88%	33.21
	$\theta_1 = 0.15$	<b>51.79%</b>	<b>44.88%</b>	<b>13.34%</b>	<b>43.70</b>
	$\theta_2 = 0$	52.55%	34.56%	34.23%	36.61
	$\theta_2 = 0.075$	52.80%	35.68%	32.42%	36.94
	$\theta_2 = 0.30$	50.26%	39.34%	21.72%	39.40

**Table 7: Variation in attack success rates, deterioration rates and average norm of attack with different sets of optimization parameters. Rows with numbers in bold represent parameter settings that achieve a good trade-off between the various performance metrics. However, both these rows have a higher average perturbation norm than the chosen set of parameters.**

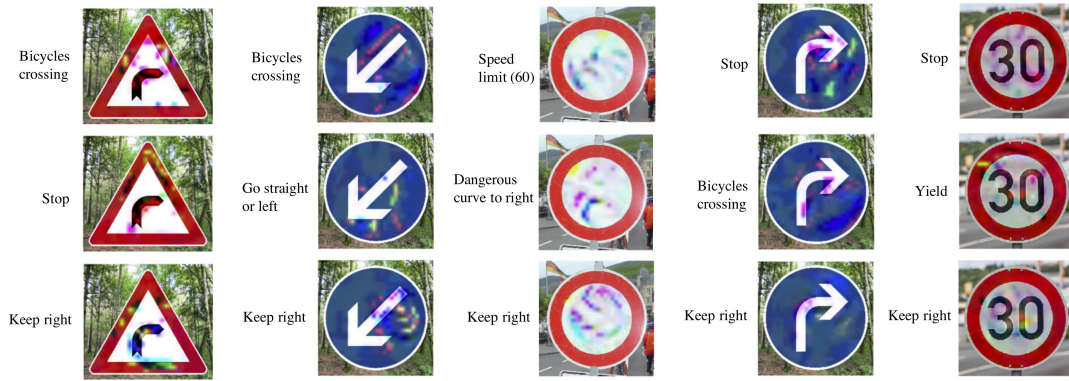
The simplest attack against DNNs is known as the Fast Gradient Sign (FGS) [14] attack which is an *untargeted attack* that involves adding a perturbation proportional to the sign of the gradient of the loss function of the neural network  $\tilde{\mathbf{x}} = \mathbf{x} + \epsilon \text{sign}(\nabla_{\mathbf{x}} \ell_f(\mathbf{x}, y))$ . The advantage of this attack is that it is extremely fast to carry out so it can be incorporated into the training phase of neural networks to make them more robust to adversarial examples.

**Note about norms:** While our attack samples have on average a maximum  $L_1$  perturbation of around 30, the  $L_1$  ball with that radius is too large to train with, since the ball has to lie in the unit hypercube. Since most of the adversarial examples do not modify each individual pixel by more than 0.3, it is a reasonable upper limit to use while training. In our experiments, we tried training with adversarial examples constrained with the  $L_1$ ,  $L_2$  and  $L_\infty$  norms and found the defense to work best with an  $L_\infty$  norm constraint.

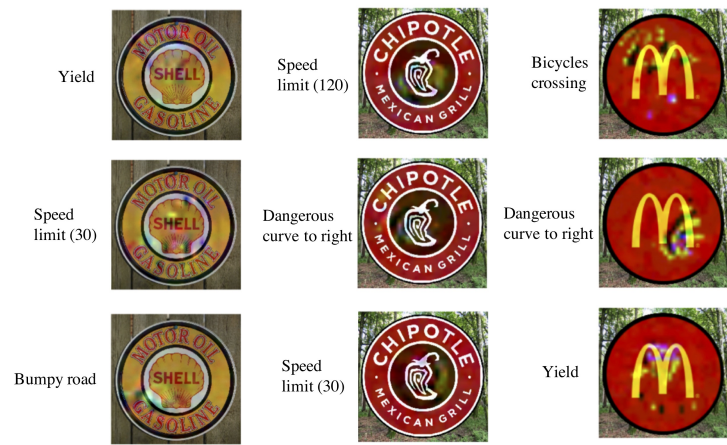
We also tried using the iterative adversarial training defense proposed by Madry et al. [51]. However, we faced a number of issues with this defense. First, it is known that the adversarial loss with iterative adversarial samples does not converge for classifiers with low capacity. When we attempted to train the Multi-scale CNN with iterative adversarial examples and the augmented training data, we observed the same behavior. Using the standard training data, the model converged to a validation accuracy of 96.2%. However, its performance on adversarial examples generated using the auxiliary

traffic data was inferior to that of the Multi-scale CNN with standard adversarial training.

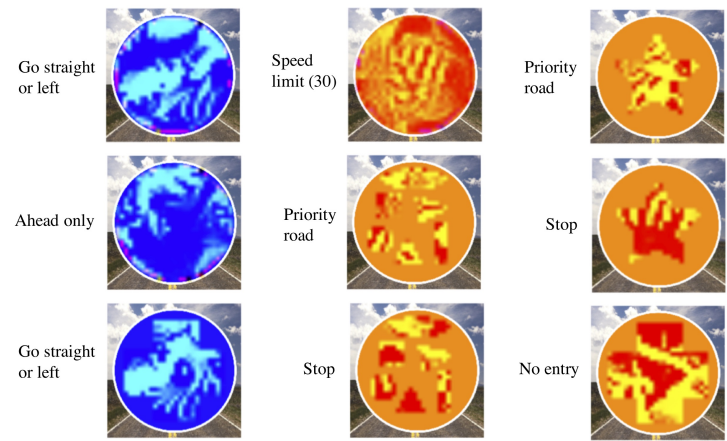




(a) Some samples of our In-Distribution (Adversarial Traffic Sign) attack, along with the labels they are classified as.



(b) Some samples of our Out-of-Distribution (Logo) attack, along with the labels they are classified as.



(c) Some samples of our Out-of-Distribution (Custom Sign) attack, along with the labels they are classified as.

Figure 14: Adversarial examples that achieve SPAS of 100%.