

---

# CURE-TSR: Challenging Unreal and Real Environments for Traffic Sign Recognition

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 In this paper, we investigate the robustness of traffic sign recognition algorithms  
2 under challenging conditions. Existing datasets are limited in terms of their size and  
3 challenging condition coverage, which motivated us to generate the Challenging  
4 Unreal and Real Environments for Traffic Sign Recognition (CURE-TSR) dataset.  
5 It includes more than two million traffic sign images that are based on real-world  
6 and simulator data. We benchmark the performance of existing solutions in real-  
7 world scenarios and analyze the performance variation with respect to challenging  
8 conditions. We show that challenging conditions can decrease the performance  
9 of baseline methods significantly, especially if these challenging conditions result  
10 in loss or misplacement of spatial information. We also investigate the utilization  
11 of simulator data along with real-world data and show that hybrid training can  
12 enhance the average recognition performance in real-world scenarios.

## 13 1 Introduction

14 Autonomous vehicles are transforming existing transportation systems. As we step up the ladder of  
15 autonomy, more critical functions are performed by algorithms, which demands more robustness.  
16 In case of following traffic rules, robust sign recognition systems are essential unless we have prior  
17 information about traffic sign types and locations. It is a common practice to test the robustness  
18 of these systems with traffic datasets (1; 2; 3; 4; 5; 6; 7; 8; 9; 10). However, majority of these  
19 datasets are limited in terms of challenging environmental conditions. There is usually no metadata  
20 corresponding to challenge conditions or levels in these datasets, which are also limited in terms of  
21 dataset size. Moreover, the relationship between challenging conditions and algorithmic performance  
22 is not analyzed in these studies. Lu *et al.* (11) investigated the traffic sign detection performance  
23 with respect to challenging adversarial examples and showed that adversarial perturbations are  
24 effective only in specific situations. Das *et al.* (12) showed the vulnerabilities of existing systems  
25 and suggested JPEG compression to eliminate adversarial effects. Even though both of these studies  
26 analyze algorithmic performance variation with respect to specific challenging situations, adversarial  
27 examples are inherently different from realistic challenging scenarios.

28 In this paper, we investigate the traffic sign recognition performance of commonly used methods under  
29 realistic challenging conditions. To eliminate the shortcomings of existing datasets, we introduce the  
30 Challenging Unreal and Real Environments for Traffic Sign Recognition (CURE-TSR) dataset. The  
31 contributions of this paper are 5 folds. First, we introduce the most comprehensive publicly-available  
32 traffic sign recognition dataset with controlled challenging conditions. Second, we provide real-world  
33 data as well as simulator data, which can enable investigating transfer learning problem between real  
34 and simulated environments. Understanding the relationship between real and simulated environments  
35 can lead to realistic dataset design, which may eventually eliminate the need for real-world data  
36 collection. Third, we provide a benchmark of commonly used methods in the introduced dataset.

37 Forth, we provide a comprehensive analysis of algorithmic performance with respect to challenging  
 38 environmental conditions. Fifth, we utilize simulator data along with real-world data to enhance the  
 39 performance of baseline methods in real-world scenarios.

## 40 2 Dataset

41 Timofte *et al.* (3) introduced the Belgium traffic sign classification (BelgiumTSC) dataset whose  
 42 images were acquired with a van that had 8 roof-mounted cameras. Acquisition vehicle cruised in  
 43 streets of Belgium and images were captured every meter. A subset of these images were selected  
 44 and traffic signs were cropped to obtain the BelgiumTSC dataset. Stallkamp *et al.* (6; 7) introduced  
 45 the German traffic sign recognition benchmark (GTSRB) dataset, which was acquired during daytime  
 46 in Germany. Each traffic sign instance in the dataset is adjusted to have 30 images. BelgiumTSC and  
 47 GTSRB datasets are limited in terms of challenging environmental conditions and they do not include  
 48 metadata related to the type of challenging conditions or their levels. Because of limited control in data  
 49 acquisition setup, it is not possible to perform controlled experiments with these datasets. The total  
 50 number of annotated signs including BelgiumTSC and GTSRB datasets is around 60,000, which may  
 51 not be sufficient to test the robustness of recognition algorithms comprehensively. To compensate  
 52 the shortcomings in the literature, we introduce the CURE-TSR dataset. Main characteristics of  
 53 BelgiumTSC, GTSRB, and CURE-TSR datasets are summarized in Table 1.

Table 1: Main characteristics of BelgiumTSC, GTSRB, and CURE-TSR datasets.

Dataset	Number of images	Number of annotated images	Number of sign types	Sign size	Origin of the videos	Acquisition device
BelgiumTSC (13)	7,095 - 7,125	All images	62	11x10 to 562x438	Captured in Belgium	Color cameras
GTSRB (14)	133,000 - 144,769	51,840	43	15x15 to 250x250	Captured in Germany	Prosilica GC 1380CH color camera
CURE-TSR	2,206,106	All images	14	3x7 to 206x277	Captured in Belgium and Generated in Unreal Engine 4	Color cameras

54



(a) Real-world (real) image



(b) Simulator (unreal) image

Figure 1: Real and unreal environments.

55

56 Traffic sign images in the CURE-TSR dataset were cropped from the CURE-TSD dataset (15), which  
 57 includes around 1.7 million real-world and simulator images. Real-world images were obtained from  
 58 the BelgiumTS video sequences and simulated images were generated with the Unreal Engine 4

59 game development tool. In Fig. 1, we show a sample real-world image and a simulator image. In the  
 60 rest of this paper, we refer to simulator generated images as unreal images and real-world images  
 61 as real images. As observed in sample images, both real and unreal images are usually from urban  
 62 environments. There are 14 traffic signs with annotations in both environments, which are shown in  
 63 Fig. 2. Sign types include speed limit, goods vehicles, no overtaking, no stopping, no parking, stop,  
 64 bicycle, hump, no left, no right, priority to, no entry, yield, and parking.



Figure 2: Traffic signs in real ( $1^{st}$  row) and unreal ( $2^{nd}$  row) environments.

65

66 Unreal and real sequences were processed with state-of-the-art visual effect software Adobe(c)  
 67 After Effects to simulate challenging conditions, which include rain, snow, haze, shadow, darkness,  
 68 brightness, blurriness, dirtiness, colorlessness, sensor and codec errors. In Fig. 3, we show sample  
 69 stop sign images under challenging conditions in both real and unreal environments.



Figure 3: Stop signs under challenging conditions in real ( $1^{st}$  row) and unreal ( $2^{nd}$  row) environments.

70

### 71 3 Experiments

#### 72 3.1 Baseline Methods, Dataset, and Performance Metric

73 In the German traffic sign recognition benchmark (GTSRB) (6), histogram of oriented gradient  
 74 (HOG) features were utilized to report the baseline results. In the Belgium traffic sign classification  
 75 (BelgiumTSC) benchmark, cropped traffic sign images were converted into grayscale and rescaled to  
 76  $28 \times 28$  patches, which were included in the baseline. Moreover, HoG features were also used as a  
 77 baseline method. They classified traffic sign images with methods including support vector machines  
 78 (SVMs). Similar to GTSRB and BelgiumTS datasets, we use rescaled grayscale and color images as  
 79 well as HoG features as baseline. In the final classification stage, we utilize one-vs-all SVMs with  
 80 radial basis kernels and softmax classifiers. In addition to aforementioned techniques, we also use a  
 81 shallow convolutional neural network, which consists of two convolutional layers followed by two  
 82 fully connected layers, and a softmax classifier. We preprocessed images using  $l_2$  normalization,  
 83 mean subtraction, and division by standard deviation.

84 Traffic sign images originate from 49 video sequences, which were split into approximately 70%  
 85 training set and 30% test set. Video sequences were split one sign at a time, starting from the least  
 86 common sign. Once video sequences were assigned to training or testing sets, splitting continued  
 87 from the remaining sequences until all the sequences were classified. In the first experiment set, we  
 88 utilized 7, 292 traffic sign images in the training stage obtained from challenge-free real training  
 89 sequences. In the testing, we utilized 3, 334 images from each challenge category and level, which  
 90 adds up to 200, 040 images ( $3, 334 \text{ images} \times 12 \text{ challenge types} \times 5 \text{ levels}$ ). As performance metric,  
 91 we utilized classification accuracy, which corresponds to the percentage of traffic signs that are  
 92 correctly classified.

93

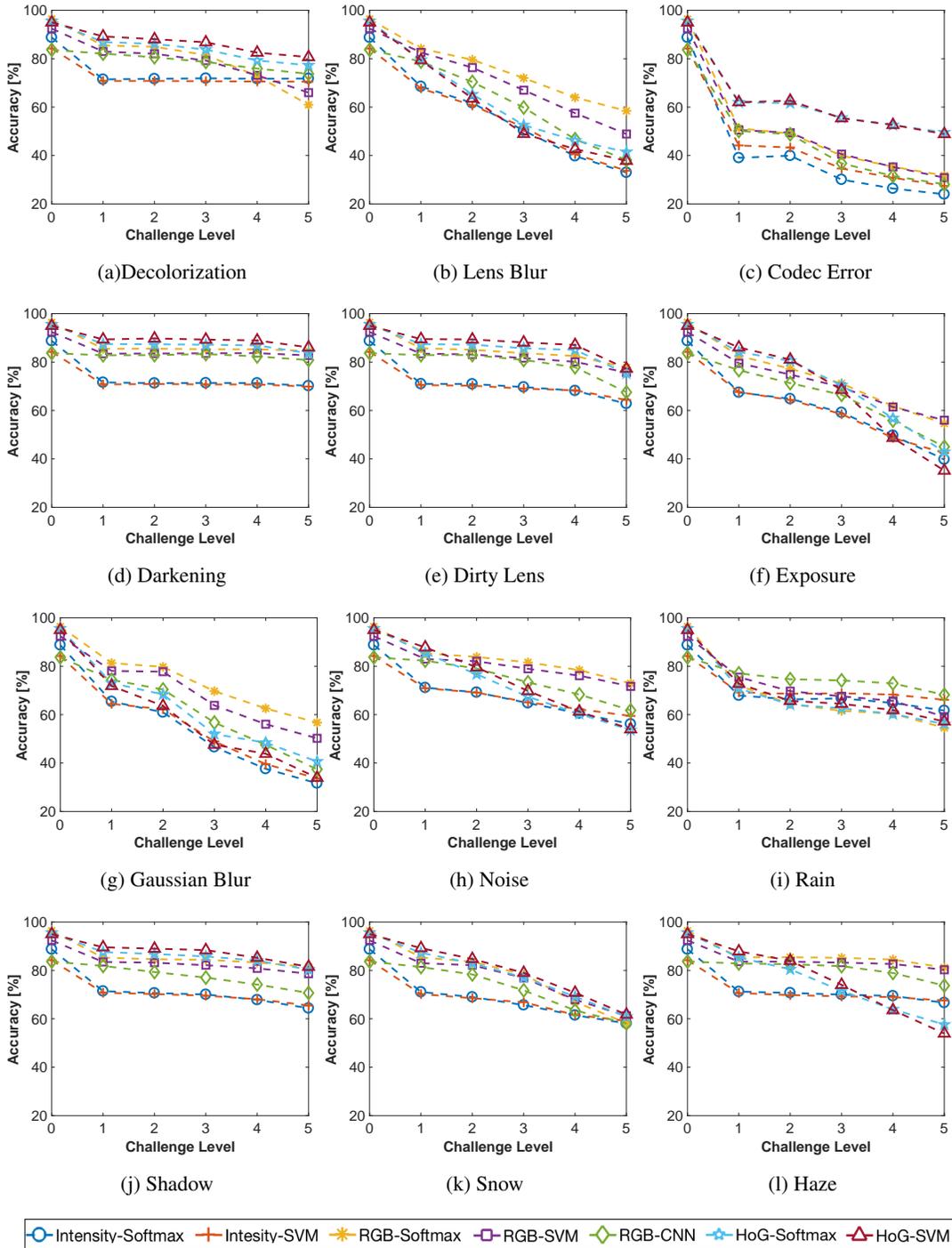


Figure 4: Performance versus challenge levels.

### 94 3.2 Experiment 1: Recognition in Real Environments under Challenging Conditions

95 We analyzed the accuracy of baseline methods with respect to challenge levels for each challenge  
96 type and report the results in Fig. 4. Severe decolorization (Fig. 4(a)) leads to at least 10% decrease  
97 in accuracy, which is less compared to majority of other challenge categories because remaining  
98 information is still sufficient for shape-based recognition. Among all the challenges, codec error is  
99 the most effective category that significantly degrades the classification accuracy even with challenge  
100 level 1 as shown in Fig. 4(c). We can observe that there is at least 30% decrease for each method after  
101 challenge level 1 and at least 46% decrease after challenge level 5. Lens blur (Fig. 4(b)), exposure  
102 (Fig. 4(f)), and Gaussian blur (Fig. 4(g)) result in significant performance decrease under severe  
103 challenging conditions, at least 36% for each baseline method. However, classification accuracy  
104 decreases more linearly in these categories compared to codec error because of its steep decrease in  
105 level 1. In darkening category (Fig. 4(d)), classification accuracy decreases at least 5% in challenge  
106 level 1 for all the methods other than CNN. The normalization operation in convolutional model makes  
107 it less sensitive to darkening challenge. When challenge level becomes more severe, performance of  
108 baseline methods degrades a few percent at most.

109 In dirty lens category (Fig. 4(e)), new dirty lens images were overlayed on entire images to increase  
110 the challenge level. And, the new dirt patterns do not necessarily occlude traffic signs. Therefore,  
111 performance of baseline methods do not always change when challenge level increases. In noise  
112 category (Fig. 4(h)), HoG and CNN correspond to a more linear performance decrease compared to  
113 intensity and color-based methods, whose performance decreases are steeper for level 1 challenge. In  
114 rain category (Fig. 4(i)), particle models are all around the scene, which result in significant occlusion  
115 even in level 1 challenge. Therefore, degradation while going from challenge-free to level 1 challenge  
116 is steeper than any further relative changes. In shadow category (Fig. 4(j)), vertical shadow lines  
117 are all over the images, which lead to relatively steep performance decrease for challenge level 1.  
118 We observe slight degradation as challenge level increases because areas under shadow become less  
119 visible. In case of snow challenge (Fig. 4(k)), intensity-based methods result in a more significant  
120 decrease compared to other methods for level 1 challenge but all methods converge to a similar  
121 classification accuracy under severe snow challenge. In haze category (Fig. 4(l)), performance of  
122 intensity-based models decrease steeply for level 1 challenge whereas decrease in HoG-based models  
123 follow a more linear behavior. Color image-based classifiers and CNN are less sensitive to haze  
124 challenge compared to other methods. Haze challenge was generated as a combination of radial  
125 gradient operator with partial opacity, a smoothing operator, an exposure operator, a brightness  
126 operator, and a contrast operator. Moreover, the location of the operator was adjusted manually per  
127 frame to simulate a sense of depth. Because of the complexity of haze model, it is less intuitive to  
128 explain the behavior of baseline methods. However, the higher tolerance of CNN model with respect  
129 to haze challenge can be explained with its capability to directly learn spatial patterns from visual  
130 representations.

### 131 3.3 Experiment 2: Recognition in Real Environments under Challenging Conditions with 132 the Help of Challenging Unreal Environments

133 In Section 3.2, we analyzed the performance of baseline methods with respect to challenging  
134 conditions and level. Baseline methods were trained with 7,292 real-world images and a total of  
135 200,040 images were used in testing. In this section, we investigate the performance of baseline  
136 methods when unreal images are used in the training in addition to real images. Test set is same  
137 as experiment 1 but we extended the training set with 20 unreal images for each traffic sign from  
138 challenge level 5 sequences. We selected the traffic signs with maximum area to obtain highest  
139 resolution samples. Overall, training set includes 3,084 unreal images (20 images  $\times$  11 challenge  
140 types  $\times$  14 traffic signs) and 7,292 real-world images. We compared the performance of baseline  
141 methods that are trained with and without unreal images and report the performance change in Table  
142 2. Each entry in the table other than the last row and the last column was obtained by calculating  
143 the performance change for a baseline method over all the challenge levels for a specific challenge  
144 type. Entries in the last row were calculated by averaging the performance change of each baseline  
145 method over all challenge types. Finally, entries in the last column were calculated by averaging the  
146 performance change over all baseline methods for each challenge type.

147

Table 2: Classification accuracy change (%) when additional unreal images used in the training.

Challenge Types	Baseline Methods						CNN	Average
	Intensity		Color		HoG			
	Softmax	SVM	Softmax	SVM	Softmax	SVM		
<b>Decolorization</b>	<b>+2.86</b>	<b>+3.32</b>	<b>+1.46</b>	-0.53	<b>+1.43</b>	-0.01	<b>+3.23</b>	<b>+1.68</b>
<b>Lens Blur</b>	<b>+3.98</b>	<b>+2.71</b>	<b>+4.45</b>	<b>+6.60</b>	<b>+3.34</b>	<b>+1.81</b>	-1.78	<b>+3.02</b>
<b>Codec Error</b>	<b>+0.47</b>	-1.21	<b>+1.51</b>	-0.82	-1.55	-1.61	<b>+2.40</b>	-0.12
<b>Darkening</b>	<b>+2.83</b>	<b>+2.98</b>	<b>+2.87</b>	<b>+1.44</b>	<b>+1.68</b>	<b>+0.44</b>	<b>+2.58</b>	<b>+2.12</b>
<b>Dirty lens</b>	<b>+3.14</b>	<b>+2.86</b>	<b>+2.68</b>	<b>+1.63</b>	<b>+2.00</b>	<b>+0.62</b>	<b>+3.11</b>	<b>+2.29</b>
<b>Exposure</b>	<b>+2.54</b>	<b>+1.77</b>	<b>+1.34</b>	<b>+1.97</b>	-0.66	-2.23	<b>+0.54</b>	<b>+0.75</b>
<b>Gaussian Blur</b>	<b>+5.89</b>	<b>+3.98</b>	<b>+4.24</b>	<b>+7.06</b>	<b>+2.03</b>	<b>+1.77</b>	<b>+2.78</b>	<b>+3.97</b>
<b>Noise</b>	<b>+1.62</b>	<b>+1.58</b>	<b>+1.89</b>	<b>+0.58</b>	<b>+1.41</b>	-0.90	<b>+2.25</b>	<b>+1.21</b>
<b>Rain</b>	<b>+2.30</b>	<b>+1.28</b>	<b>+4.73</b>	<b>+2.75</b>	<b>+5.48</b>	<b>+2.34</b>	<b>+0.69</b>	<b>+2.80</b>
<b>Shadow</b>	<b>+2.95</b>	<b>+3.38</b>	<b>+3.27</b>	<b>+1.62</b>	<b>+1.73</b>	<b>+0.64</b>	<b>+3.01</b>	<b>+2.37</b>
<b>Snow</b>	<b>+3.19</b>	<b>+2.81</b>	<b>+2.09</b>	<b>+0.48</b>	<b>+2.63</b>	<b>+0.92</b>	<b>+4.34</b>	<b>+2.35</b>
<b>Haze</b>	<b>+3.28</b>	<b>+3.22</b>	<b>+3.22</b>	<b>+1.41</b>	<b>+2.26</b>	-1.35	<b>+3.51</b>	<b>+2.22</b>
<b>All (average)</b>	<b>+2.92</b>	<b>+2.39</b>	<b>+2.81</b>	<b>+2.02</b>	<b>+1.81</b>	<b>+0.20</b>	<b>+2.22</b>	-

148 We tested 7 baseline methods over 12 challenge types and report the performance change of each  
 149 baseline method for each challenge type. Out of 84 result categories (7 baseline methods  $\times$  12  
 150 challenge types), classification performance increased in 72 of them. On average, classification  
 151 performance increased for all challenge types other than a slight decrease in codec error. Moreover,  
 152 average classification performance increased for each baseline method, which is a slight increase  
 153 for HoG-SVM (0.2%) and more for other methods (at least 1.81%). Additional unreal images  
 154 in the training set were obtained from all the challenge types except haze category. However,  
 155 classification accuracy increased for all the baseline methods at least 1.41% other than HoG-SVM  
 156 in haze category. The performance enhancement in haze can be understood by analyzing the  
 157 computational model of haze and its perceptual similarity to other challenges. Haze model includes  
 158 a smoothing operator, an exposure filter, a brightness operator, and a contrast operator. Exposure  
 159 filter is used in the exposure (overexposure) model and smoothing operator is utilized in blur models.  
 160 Moreover, perceptually, we can observe similarities between haze and blur challenges in terms of  
 161 smoothness and similarities between haze and exposure in terms of washed out details. Therefore,  
 162 perceptually and computationally similar challenges in the training stage can affect the performance  
 163 of each other in the testing stage.

## 164 4 Conclusion

165 We introduced the CURE-TSR dataset, which is the most comprehensive traffic sign recognition  
 166 dataset in the literature that includes controlled challenging conditions. We provided a benchmark of  
 167 commonly used methods in the CURE-TSR dataset and reported that challenging conditions leads to  
 168 severe performance degradation for all baseline methods. We have shown that lens blur, exposure,  
 169 Gaussian blur, and codec error degrade recognition performance more significantly compared to  
 170 other challenge types because these challenge categories directly result in losing or misplacing  
 171 shape-related information. In addition to training and testing data-driven methods with real-world  
 172 data, we also utilized simulator images in the training and reported performance enhancement for  
 173 most of the baseline methods and challenge categories in real-world scenarios.

## References

- 174
- 175 [1] C. Grigorescu and N. Petkov, "Distance sets for shape filters and shape recognition," *IEEE*  
176 *Trans. Image Proces.*, vol. 12, no. 10, pp. 1274–1286, Oct 2003.
- 177 [2] R. Timofte, K. Zimmermann, and L. V. Gool, "Multi-view traffic sign detection, recognition,  
178 and 3D localisation," in *WACV*, Dec 2009, pp. 1–8.
- 179 [3] R. Timofte, K. Zimmermann, and L. Van Gool, "Multi-view traffic sign detection, recognition,  
180 and 3D localisation," *Mach. Vis. App.*, vol. 25, no. 3, pp. 633–647, 2014.
- 181 [4] R. Belaroussi, P. Foucher, J. P. Tarel, B. Soheilian, P. Charbonnier, and N. Paparoditis, "Road  
182 sign detection in images: A case study," in *Proc. ICPR*, Aug 2010, pp. 484–488.
- 183 [5] F. Larsson and M. Felsberg, "Using fourier descriptors and spatial models for traffic sign  
184 recognition," in *Proc. SCIA*, Berlin, Heidelberg, 2011, SCIA'11, pp. 238–249, Springer-Verlag.
- 185 [6] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "The German Traffic Sign Recognition  
186 Benchmark: A multi-class classification competition," in *Proc. IJCNN*, July 2011, pp. 1453–  
187 1460.
- 188 [7] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "Man versus computer: Benchmarking  
189 machine learning algorithms for traffic sign recognition," *Neural Networks*, vol. 32, pp. 323 –  
190 332, 2012, Selected Papers from IJCNN 2011.
- 191 [8] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel, "Detection of traffic signs in  
192 real-world images: The German Traffic Sign Detection Benchmark," in *Proc. IJCNN*, Aug  
193 2013, pp. 1–8.
- 194 [9] A. Mogelmoose, M. M. Trivedi, and T. B. Moeslund, "Vision-based traffic sign detection and  
195 analysis for intelligent driver assistance systems: Perspectives and survey," *IEEE Trans. Intell.*  
196 *Transp. Syst.*, vol. 13, no. 4, pp. 1484–1497, Dec 2012.
- 197 [10] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu, "Traffic-sign detection and classification  
198 in the wild," in *Proc. IEEE CVPR*, June 2016, pp. 2110–2118.
- 199 [11] J. Lu, H. Sibai, E. Fabry, and D. Forsyth, "NO need to worry about adversarial examples in  
200 object detection in autonomous vehicles," in *arXiv:1707.03501*, 2017.
- 201 [12] N. Das, M. Shanbhogue, S.-T. Chen, F. Hohman, L. Chen, M. E. Kounavis, and D. H. Chau,  
202 "Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression," in  
203 *arXiv:1705.02900*, 2017.
- 204 [13] R. Timofte, K. Zimmermann, and L. V. Gool, "Belgium traffic sign dataset," [http://btsd.  
205 ethz.ch/shareddata/](http://btsd.ethz.ch/shareddata/).
- 206 [14] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "German traffic sign recognition and  
207 detection benchmarks," <http://benchmark.ini.rub.de/>.
- 208 [15] "CURE-TSD: Challenging unreal and real environments for traffic sign detection," [https:  
209 //ghassanalregib.com/cure-tsd/](https://ghassanalregib.com/cure-tsd/).