

**Citation** D. Temel, G. Kwon\*, M. Prabhushankar\*, and G. AlRegib, “CURE-TSR: Challenging unreal and real environments for traffic sign recognition,” in Neural Information Processing Systems (NIPS) Workshop on Machine Learning for Intelligent Transportation Systems (MLITS), Long Beach, U.S., December 2017.

**Review** Date of acceptance to NIPS MLITS : November 14 2017

**Data/Codes** <https://ghassanalregib.com/cure-tsr/>

**Bib** @INPROCEEDINGS{Temel2017\_NIPS,  
author={D. Temel and G. Kwon and M. Prabhushankar and G. AlRegib},  
booktitle={Neural Information Processing Systems (NIPS) Workshop on Machine Learning for Intelligent Transportation Systems (MLITS)},  
title={CURE-TSR: Challenging unreal and real environments for traffic sign recognition},  
year={2017},  
month={December}, }

**Contact** [alregib@gatech.edu](mailto:alregib@gatech.edu) <https://ghassanalregib.com/>  
[dcantemel@gmail.com](mailto:dcantemel@gmail.com) <http://cantemel.com/>

---

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

---

Dogancan Temel, Gukyeong Kwon\*, Mohit Prabhushankar\*, and Ghassan AlRegib  
Georgia Institute of Technology  
Center for Signal and Information Processing  
{cantemel, gukyeong.kwon, mohit.p, alregib}@gatech.edu

## Abstract

In this paper, we investigate the robustness of traffic sign recognition algorithms under challenging conditions. Existing datasets are limited in terms of their size and challenging condition coverage, which motivated us to generate the Challenging Unreal and Real Environments for Traffic Sign Recognition (**CURE-TSR**) dataset. It includes more than two million traffic sign images that are based on real-world and simulator data. We benchmark the performance of existing solutions in real-world scenarios and analyze the performance variation with respect to challenging conditions. We show that challenging conditions can decrease the performance of baseline methods significantly, especially if these challenging conditions result in loss or misplacement of spatial information. We also investigate the effect of data augmentation and show that utilization of simulator data along with real-world data enhance the average recognition performance in real-world scenarios. The dataset is publicly available at <https://ghassanalregib.com/cure-tsr/>.

## 1 Introduction

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

In this paper, we investigate the traffic sign recognition performance of commonly used methods under realistic challenging conditions. To eliminate the shortcomings of existing datasets, we introduce the Challenging Unreal and Real Environments for Traffic Sign Recognition (CURE-TSR) dataset. The contributions of this paper are 4 folds.

- We introduce the most comprehensive publicly-available traffic sign recognition dataset with controlled challenging conditions.

---

\*Equal contribution.

- We provide a detailed analysis of the benchmarked algorithms in terms of their recognition performance under challenging conditions. Based on this analysis, we identify the vulnerabilities of algorithms with respect to challenging conditions, which should give insights into the use of such models under certain conditions.
- We provide images that originate from captured sequences as well as synthesized sequences, which would lead to a better understanding of similarities/differences between real-world and simulator data in terms of algorithmic performance. This understanding can be utilized to generate more realistic datasets and minimize the need for real-world data collection that requires significant resources.
- We use diverse data augmentation methods and show that utilization of limited simulator data along with real-world data can enhance the recognition performance. This observation shows that simulated environments can enhance the performance of data-driven methods in real-world scenarios even when there is a difference between target and source domains.

## 2 Dataset

Timofte *et al.* [3] introduced the Belgium traffic sign classification (BelgiumTSC) dataset whose images were acquired with a van that had 8 roof-mounted cameras. Acquisition vehicle cruised in streets of Belgium and images were captured every meter. A subset of these images were selected and traffic signs were cropped to obtain the BelgiumTSC dataset. Stallkamp *et al.* [6, 7] introduced the German traffic sign recognition benchmark (GTSRB) dataset, which was acquired during daytime in Germany. Each traffic sign instance in the dataset is adjusted to have 30 images. BelgiumTSC and GTSRB datasets are limited in terms of challenging environmental conditions and they do not include metadata related to the type of challenging conditions or their levels. Because of limited control in data acquisition setup, it is not possible to perform controlled experiments with these datasets. The total number of annotated signs including BelgiumTSC and GTSRB datasets is around 60,000, which may not be sufficient to test the robustness of recognition algorithms comprehensively. To compensate the shortcomings in the literature, we introduce the CURE-TSR dataset. Main characteristics of 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 [15]	2,206,106	All images	14	3x7 to 206x277	Captured in Belgium and Generated in Unreal Engine 4	Color cameras

Traffic sign images in the CURE-TSR dataset were cropped from the CURE-TSD dataset [16], which includes around 1.7 million real-world and simulator images with more than 2 million traffic sign instances. Real-world images were obtained from the BelgiumTS video sequences and simulated images were generated with the Unreal Engine 4 game development tool. In Fig. 1, we show a sample real-world image and a simulator image. In the rest of this paper, we refer to simulator generated images as unreal images and real-world images as real images. As observed in sample images, both real and unreal images are usually from urban environments. While deciding on the type of traffic signs to be included in real and unreal sequences, we focused on two main criteria. First, not every sign type can be reasonably located in unreal sequences. Second, there are limited number of common signs between the package utilized in the generation of unreal sequences and real sequences. Based on the aforementioned selection criteria, we narrowed down number of traffic signs to 14 types as shown in Fig. 2. Sign types include *speed limit*, *goods vehicles*, *no overtaking*, *no stopping*, *no parking*, *stop*, *bicycle*, *hump*, *no left*, *no right*, *priority to*, *no entry*, *yield*, and *parking*.



(a) Real-world (real) image



(b) Simulator (unreal) image

Figure 1: Real and unreal environments.



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

Unreal and real sequences were processed with state-of-the-art visual effect software Adobe(c) After Effects to simulate challenging conditions, which include rain, snow, haze, shadow, darkness, brightness, blurriness, dirtiness, colorlessness, sensor and codec errors. The key component in this study is not the number of traffic signs but the number of challenging conditions and the context of each traffic sign in a virtual dataset and its corresponding real dataset. If one considers a traffic sign in a challenging condition as a distinct configuration, then we end up with 182 ( $14 \times 13$ ) distinct configurations in real sequences and 168 ( $14 \times 12$ ) distinct configurations in virtual sequences. In Fig. 3, we show sample stop sign images under challenging conditions in both real and unreal environments. We included `codec error` as an edge case to test the limits of benchmarked methods. Recognizing traffic signs with `codec errors` can be challenging even for subjects because of significant misalignment. If a sign is totally misaligned, it will not be possible to recognize it at cropped location but in case there is residual, it can still be possible to recognize that traffic sign. `Codec-related errors` can be critical in various applications including but not limited to remote driving. Overall, there are 5 challenge levels for each challenge category, which are shown in Appendix A.



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

### 3 Experiments

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

In the German traffic sign recognition benchmark (GTSRB) [6], histogram of oriented gradient (HOG) features were utilized to report the baseline results. In the Belgium traffic sign classification (BelgiumTSC) benchmark, cropped traffic sign images were converted into grayscale and rescaled to

$28 \times 28$  patches, which were included in the baseline. Moreover, HoG features were also used as a baseline method. They classified traffic sign images with methods including support vector machines (SVMs). Similar to GTSRB and BelgiumTSC datasets, we use rescaled grayscale and color images as well as HoG features as baseline. In the final classification stage, we utilize one-vs-all SVMs with radial basis kernels and softmax classifiers. In addition to aforementioned techniques, we also use a shallow convolutional neural network, which consists of two convolutional layers followed by two fully connected layers, and a softmax classifier. We preprocessed images using  $l_2$  normalization, mean subtraction, and division by standard deviation.

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

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

We analyze the accuracy of baseline methods with respect to challenge levels for each challenge type and report the results in Fig. 4. Severe decolorization (Fig. 4(a)) leads to at least 10% decrease in accuracy for color-based and HoG-based methods. However, intensity-based methods show consistent performance over different challenge levels since no color information is used by intensity-based methods. Among all the challenges, `codec_error` is the most effective category that significantly degrades the classification accuracy even with challenge level 1 as shown in Fig. 4(c). We can observe that there is at least 30% decrease for each method after challenge level 1 and at least 46% decrease after challenge level 5. `lens_blur` (Fig. 4(b)), `exposure` (Fig. 4(f)), and `Gaussian blur` (Fig. 4(g)) result in significant performance decrease under severe challenging conditions, at least 27% for each baseline method. However, classification accuracy decreases more linearly in these categories compared to `codec_error` because of its steep decrease in level 1. In `darkening` category (Fig. 4(d)), classification accuracy is consistent for all the methods. The normalization operation in the preprocessing step makes all methods less sensitive to `darkening` challenge. When challenge level becomes more severe, performance of baseline methods degrades a few percent at most.

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

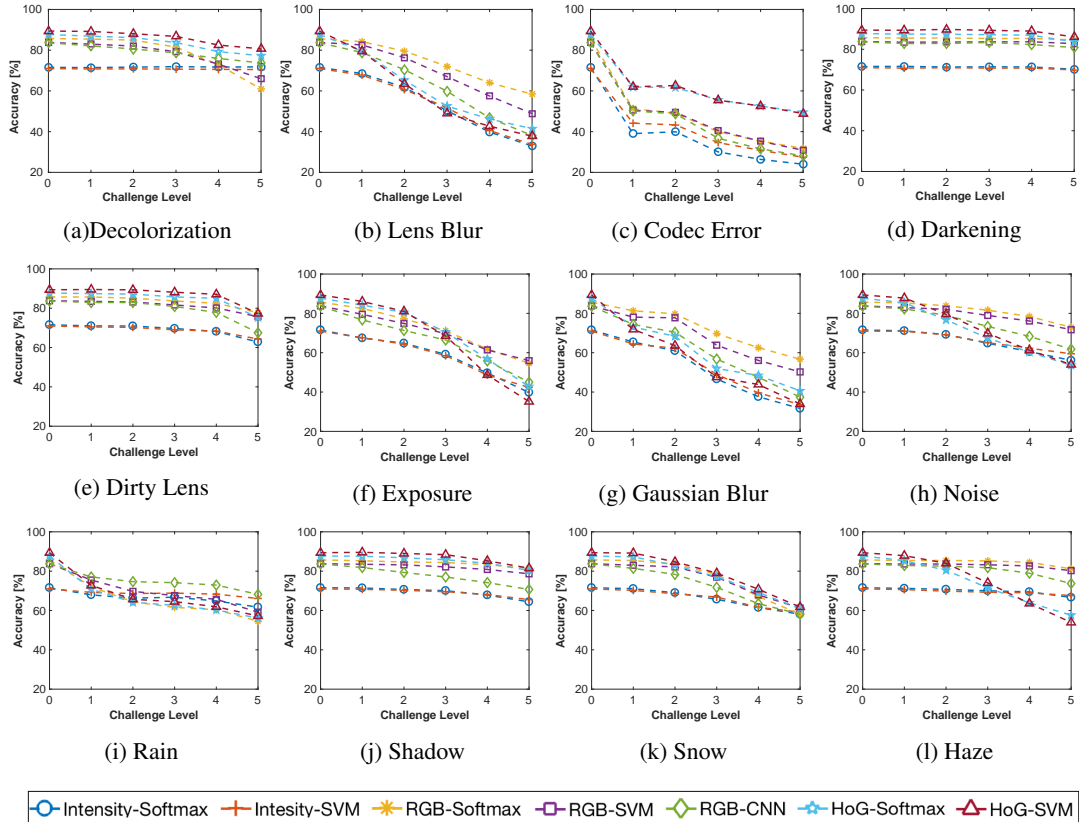


Figure 4: Performance versus challenge levels.

### 3.3 Experiment 2: Recognition in Real Environments under Challenging Conditions with Data Augmentation

We investigate the role of data augmentation methods in traffic sign recognition under challenging conditions. Augmented data include flipped real images, real challenge images, and unreal challenge images. To augment flipped images, real traffic sign images were randomly flipped horizontally, vertically or horizontally and vertically. To augment real challenge images, we selected 20 traffic sign images with maximum area (highest resolution samples) for each traffic sign in the training set. Then, we obtained corresponding images for each challenge type and level. It should be noted that augmented data is challenge-free version of the challenge-free data, which is already utilized in the training. We utilized the same source images because challenging conditions in the original video sequences were synthesized globally over entire videos and it is not possible to apply the same challenge generation framework directly over new traffic sign images. Overall, in both augmentation experiments, training set include 3,080 images (20 images  $\times$  11 challenge types  $\times$  14 traffic signs) and 7,292 (complete challenge-free training set) real images. Test set is same as experiment 1 for all three data augmentation methods.

Data augmentation based on flipping degrades the recognition performance by 6.78%. Decrease in recognition performance can be mainly because of asymmetric characteristics of traffic sign images. We further explain details about this performance decrease in Appendix B. Data augmentation with real challenge images slightly decreases the average performance by 0.45%. Even though novel challenging conditions are added in the augmentation stage, original images are already included in the training. Therefore, such augmentation method does not lead to any performance enhancement in tested scenarios. Similar to real challenge images, we obtained unreal challenge images by selecting the traffic signs with maximum area for each traffic sign, challenge type, and challenge level. Instead of using solely 20 distinct images for each sign, we utilize 220 distinct unreal images (one distinct image for each sign and challenge category) in the data augmentation, which enriches training set with new angle, contrast, and lighting configurations. Results of unreal image-based data augmentations



are summarized in Table 2. Each entry in the table other than the last row and the last column was obtained by calculating the performance change for a baseline method over all the challenge levels for a specific challenge type. Entries in the last row were calculated by averaging the performance change of each baseline method over all challenge types. Finally, entries in the last column were calculated by averaging the performance change over all baseline methods for each challenge type.

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>	-

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

## 4 Conclusion

We introduced the CURE-TSR dataset, which is the most comprehensive traffic sign recognition dataset in the literature that includes controlled challenging conditions. We provided a benchmark of commonly used methods in the CURE-TSR dataset and reported that challenging conditions lead to severe performance degradation for all baseline methods. We have shown that `lens blur`, `exposure`, `Gaussian blur`, and `codec error` degrade recognition performance more significantly compared to other challenge types because these challenge categories directly result in losing or misplacing structural information. We also investigated the effect of data augmentation and showed that flipping or simply adding challenging conditions to training data do not necessarily enhance recognition performance. However, experimental results showed that utilization of diverse simulator data with challenging conditions can enhance the average recognition performance in real-world scenarios.

## References

- [1] C. Grigorescu and N. Petkov. Distance sets for shape filters and shape recognition. *IEEE Trans. Image Proces.*, 12(10):1274–1286, Oct 2003.
- [2] R. Timofte, K. Zimmermann, and L. V. Gool. Multi-view traffic sign detection, recognition, and 3D localisation. In *WACV*, pages 1–8, Dec 2009.
- [3] R. Timofte, K. Zimmermann, and L. Van Gool. Multi-view traffic sign detection, recognition, and 3D localisation. *Mach. Vis. App.*, 25(3):633–647, 2014.
- [4] R. Belaroussi, P. Foucher, J. P. Tarel, B. Soheilian, P. Charbonnier, and N. Paparoditis. Road sign detection in images: A case study. In *Proc. ICPR*, pages 484–488, Aug 2010.
- [5] F. Larsson and M. Felsberg. Using fourier descriptors and spatial models for traffic sign recognition. In *Proc. SCIA, SCIA’11*, pages 238–249, Berlin, Heidelberg, 2011. Springer-Verlag.
- [6] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. The German Traffic Sign Recognition Benchmark: A multi-class classification competition. In *Proc. IJCNN*, pages 1453–1460, July 2011.
- [7] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. Man versus computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural Networks*, 32:323 – 332, 2012. Selected Papers from IJCNN 2011.
- [8] S. Houben, J. Stallkamp, J. Salmen, M. Schlipsing, and C. Igel. Detection of traffic signs in real-world images: The German Traffic Sign Detection Benchmark. In *Proc. IJCNN*, pages 1–8, Aug 2013.
- [9] 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 Trans. Intell. Transp. Syst.*, 13(4):1484–1497, Dec 2012.
- [10] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu. Traffic-sign detection and classification in the wild. In *Proc. IEEE CVPR*, pages 2110–2118, June 2016.
- [11] J. Lu, H. Sibai, E. Fabry, and D. Forsyth. No need to worry about adversarial examples in object detection in autonomous vehicles. In *arXiv:1707.03501*, 2017.
- [12] N. Das, M. Shanbhogue, S.-T. Chen, F. Hohman, L. Chen, M. E. Kounavis, and D. H. Chau. Keeping the bad guys out: Protecting and vaccinating deep learning with jpeg compression. In *arXiv:1705.02900*, 2017.
- [13] R. Timofte, K. Zimmermann, and L. V. Gool. Belgium traffic sign dataset. <http://btsd.ethz.ch/shareddata/>, 2009.
- [14] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel. German traffic sign recognition and detection benchmarks. <http://benchmark.ini.rub.de/>, 2013.
- [15] CURE-TSR: Challenging unreal and real environments for traffic sign recognition. <https://ghassanalregib.com/cure-tsr/>, 2017.
- [16] CURE-TSD: Challenging unreal and real environments for traffic sign detection. <https://ghassanalregib.com/cure-tsd/>, 2017.



## A Appendix: Visualization of Challenge Levels and Types

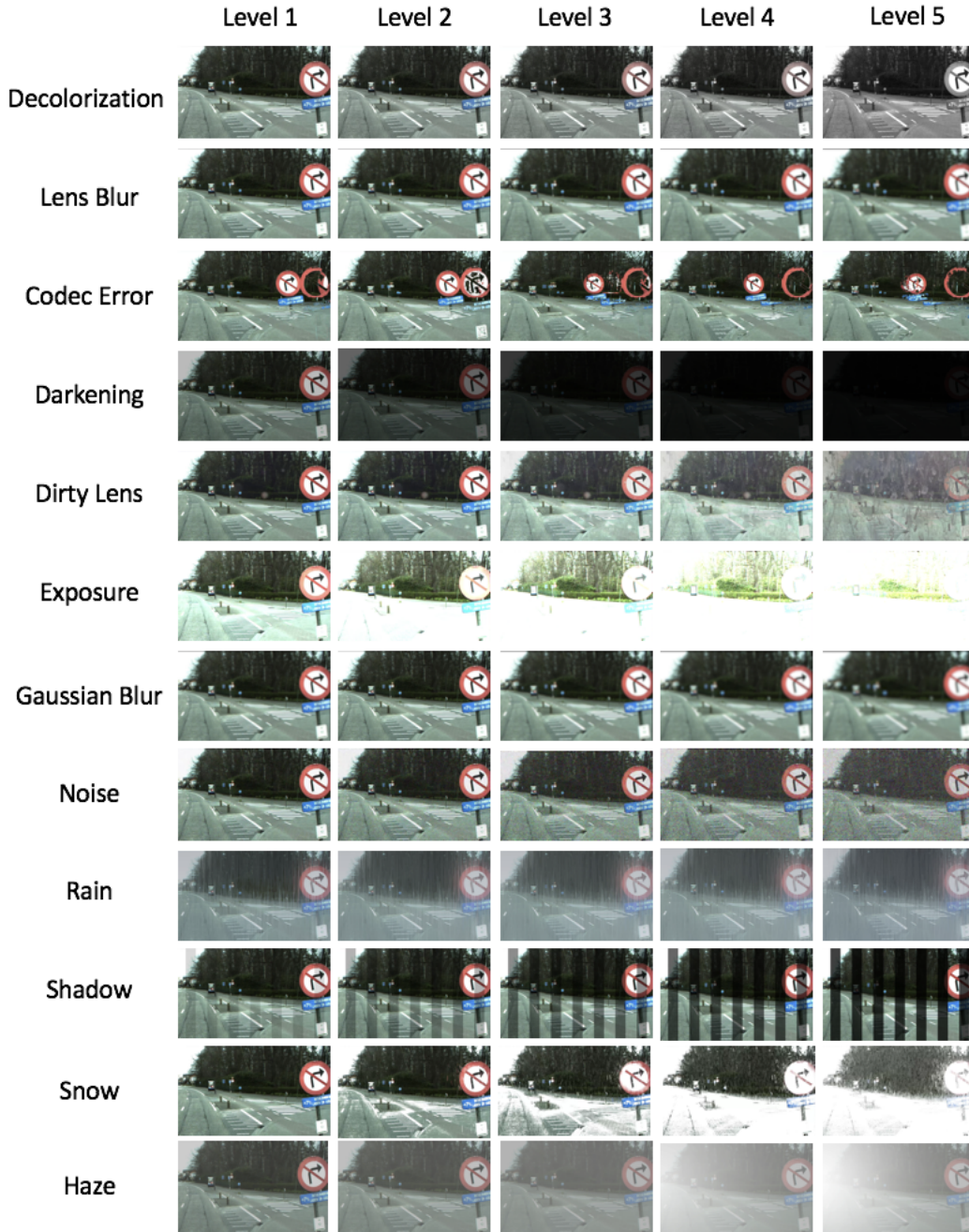


Figure 5: Challenging scene examples from each challenge type and level in CURE-TSD [16] and CURE-TSR [15] datasets.

To visualize scenes with realistic challenging condition types and levels, we cropped surrounding environments with traffic signs as shown in Fig. 5. Each row corresponds to a challenging condition and each column corresponds to a certain level of the challenging condition. Compared to other existing datasets, the CURE-TSR dataset contains more diverse challenging conditions and levels, which enables a comprehensive platform to test the robustness of recognition algorithms under challenging conditions.

## B Appendix: Data Augmentation

We retrained the benchmark algorithms, listed in Sec. 3.1, by augmenting the initial training images with their flipped versions. Vertically, horizontally or vertically and horizontally flipped challenge-free real images were used for data augmentation. Translation was not utilized in the data augmentation because recognition dataset is based on cropped images that do not include background information. Flipping-based data augmentation degrades the average recognition accuracy by more than 6.5% mainly because of the asymmetric characteristics of traffic sign images. For instance, consider the *no stopping* and *no parking* signs from Fig. 2. The former sign is symmetric along its horizontal axis while the latter sign is asymmetric. Augmenting the training data with horizontally flipped versions of the asymmetric *no parking* can lead to learning a visual representation similar to *no stopping* sign, which is different from the intended class and can degrade the overall degradation accuracy.

We visualize two softmax RGB trained models, one of which was trained without data augmentation while the other was trained with data augmentation (flipped real images) in Figs. 6 and 7 respectively. Consider the case of the learned *no parking* model (2<sup>nd</sup> row, 1<sup>st</sup> column, green highlight). Perceptually, the data-augmented model has two diagonal lines crossing each other as opposed to the single diagonal in the model learned without data augmentation. However, lines that perpendicularly cross is a characteristic of the learned *no stopping* model (1<sup>st</sup> row, 4<sup>th</sup> column, yellow highlight), which would result in misclassification. Increase in the misclassification rate between these two signs because of flipping-based data augmentation can be understood from the confusion matrices in Fig. 6(b) and Fig. 7(b) (highlighted in yellow: class types 4 and 5) in which darker colors correspond to more misclassification.

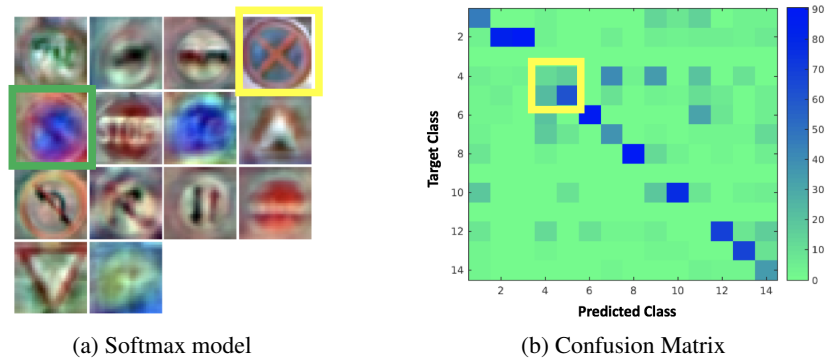


Figure 6: RGB softmax model visualization and averaged confusion matrix without data augmentation.

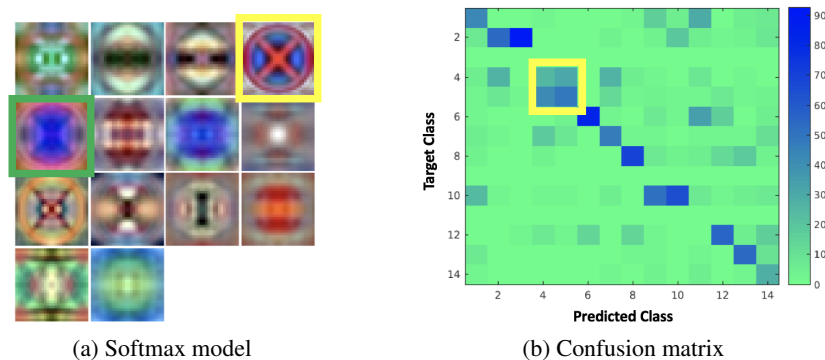


Figure 7: RGB softmax model visualization and averaged confusion matrix with data augmentation.