

# Deep Learning for Automotive Sign Detection

B.Karthiga<sup>1</sup>, M.Sivakumar<sup>2</sup>, S.Sivasubramanian<sup>3</sup>, S.Venkatesan<sup>4</sup>, V.Baranidharan<sup>5</sup>

Department of Electronics and Communication Engineering,  
 Dhanalakshmi Srinivasan Engineering College (Autonomous), Perambalur, India.

Email: <sup>1</sup>karthiga.jaya15@gmail.com, <sup>2</sup>esivakumarm@gmail.com,

<sup>3</sup>ssivasubramanian@gmail.com, <sup>4</sup>svenkatesanv7@gmail.com, <sup>5</sup>jacksbani303@gmail.com

**Abstract**—Traffic sign and light detections are core components of Advanced Driver Assistance Systems (ADAS) and self-driving vehicles. To this end, the automotive industry is widely exploiting computer vision (CV) and deep learning (DL) techniques. This paper presents a lightweight traffic sign and light detector by harnessing a single-stage, single-shot multi-object detector (SSD). For accelerating the inference speed of the detector, its original backbone, VGG16 is replaced by MobileNet V2 that expertly manages detection speed and network size. In autonomous driving, quicker detection performance with respect to the distance of an object is of particular interest, for a comfortable braking. However, farther distance makes the objects to be detected appear smaller. Unfortunately, the original SSD struggles to detect small objects. Thus, this work further optimizes the number of feature map layers of the SSD for the detection of small objects along with a better trade-off between detection precision and inference time. Experimental analysis confirms the effectiveness of the proposed model, which achieves 2 times (or more) faster detection time than the baseline SSD models and a competitive precision of 76.7%.

**Index Terms**—Computer vision, object detection, ADAS, deep learning

## I. INTRODUCTION

A reliable real-time detection of traffic sign and light on computationally limited platforms is an important concern for autonomous driving. Therefore, only the models with fewer parameters and low computational complexity are needed. In this line, various deep learning frameworks, two-stage and single-stage object detectors have been proposed by the research community. For example, the Darknet [2] is a simple architecture, having fast inference speed. But it specifically supports only NVIDIA CUDA for acceleration. Hence, the two-stage models, like the faster region-based convolutional neural networks (R-CNN) [3] and region-based fully convolutional networks (R-FCN) [4] have better detection rate; however, they require immense computational power that makes them unacceptable for real-time applications. On the other hand, the single-stage models bridge region proposal, classification and regression tasks, as a single multi-task learning. For instance, the you only look once (YOLO) [5] and single-shot multi-box detector (SSD) [1]. These models improve the detection speed and therefore can be implemented on embedded platforms. However, their efficiency is lower in terms of accuracy when compared with two-stage detectors.

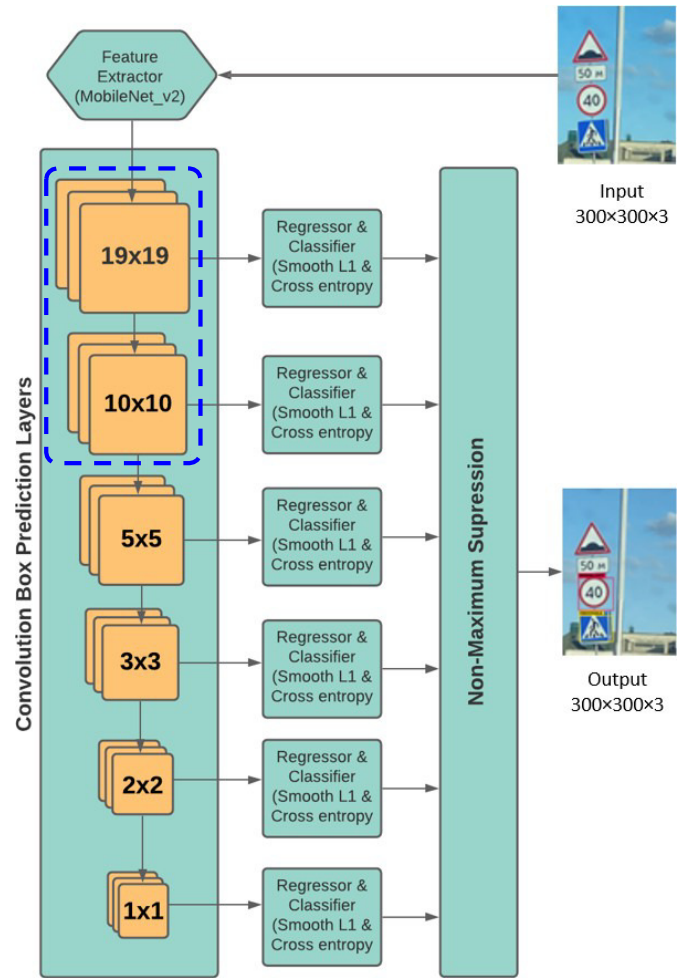


Fig. 1. The proposed SSD-based traffic sign and light detection framework. The blue broken lines indicate the only two feature map layers exploited by the proposed model out of the 6-box prediction layers in the original SSD model in [1]. The rationale of this modification is investigated in Section III.

There is another major problem with SSD that it is not good at dealing with small objects. Regardless of their small size, the traffic signs and lights provide key contextual information for intelligent transportation and safety. This work addresses these issues of SSD. The proposed detection framework integrates the multi-scale feature maps obtained by a backbone network

to enhance the detection performance. The original SSD uses six feature map layers of different scales; however, it struggles in small object detection when all six feature map layers are used. Through exhaustive experimental investigation we found that the first two feature map layers ( $19 \times 19$  and  $10 \times 10$ , cf. Fig. 1) exclusively provides not only better results in terms of small object detection but also significant reduction in detection time. To further improve the inference speed, the base network of SSD that is originally the VGG16 [6] is also replaced by MobileNetv2 [7] architecture.

The rest of the paper is organized as follows. Section II briefly describes the state-of-the-art models. Section III elaborates the proposed SSD-based approach and the improvements made in this work. Section IV presents a thorough experimental analysis. Section V closes the paper with conclusion and future direction.

## II. RELATED WORK

Most of the existing solutions focus on a single objective either the detection of the traffic sign or traffic light; not on both. In retrospect, each of the state-of-the-art for traffic sign and traffic light detection is separated by two approaches: conventional models and deep learning (DL) models.

### A. Conventional Models

These methods solely depend on the feature extraction algorithms. Let it be a classification or detection problem, they mandatorily require low-level features, like color intensities, edge details, and shape features. For example, Kim *et al.* [8] and Diaz-Cabrera *et al.* [9] devise a traffic light detection system using color information, viz. RGB, hue-saturation-intensity (HSI), and hue-saturation-value (HSV). On the other hand, the researchers in [10], handle the traffic light detection as a shape extraction of the traffic light. Similarly, Swathi *et al.* [11] and Supreesh *et al.* [12] also used color clues, as the main features to locate the traffic signs on the road scenes. Hence, Nguyen *et al.* [13] and Yang *et al.* [14] take advantage of shape descriptors, like Hough transform to extract class-specific features, such as circles, and rectangles. Later they use these features to train a machine learning model to locate the traffic signs on a given image. These models are feature-specific, and although they achieve a higher precision, they lack robustness and do not generalize well as a whole system.

### B. Deep Learning Models

The advanced object detection, segmentation, and classification algorithms for intelligent transportation application exploits deep learning [15]–[17]. Thus, researchers have focussed on building deep neural networks (DNNs), for traffic light and sign detection. For instance, Weber *et al.* [18] implement model, coined as DeepTLR. The network returns a pixel-segmented image and then they apply a bounding box regressor for detection of traffic lights. Similarly, Behrendt *et al.* [19] introduce a system for detection, tracking, and classification using a deep convolutional neural network (DCNN). By utilizing the general object detector, the SSD, Muller *et al.* [20]

introduce a model only for traffic light detection. Following the work of [19], Yudin *et al.* [21] propose another traffic light detector based on fully convolutional network (FCN). They use the FCN to get a heat-map highlighting plausible areas of traffic lights, then employ a high-speed clustering algorithm to obtain traffic light bounding boxes. Besides being a transfer learning approach, it has a very low precision of detection as compared to SSD-based solutions, like in [20].

On the other hand, for traffic sign detection, Zhang *et al.* [22] use a modified version of YOLOv2 [23] object detector. They manipulate the size of filters, and the number of layers to obtain a balance between detection accuracy and speed. Zhu *et al.* [24] design a custom-built CNN architecture to target more on the smaller size objects. Such target-specific implementation also faces lack of generalization performance. These solutions can achieve good robustness compared to the conventional counterparts, as they self learn the feature correspondence between the raw inputs and targets. However, they face the criticism for being hungry for data and compute power.

## III. PROPOSED TRAFFIC LIGHT AND SIGN DETECTION FRAMEWORK

### A. Basic Concept

For object detection, a CNN-based feature extractor (cf. Fig. 1) is extended to a larger network by eradicating top classification layers of the base network and adding some successive layers. The successive layers are connected to two main heads: (1) a regressor to predict bounding boxes, (2) a classifier to classify each of the detected boxes. Then, a non-maximum suppression (NMS) algorithm is applied to discard insignificant regions, finalizing the most probable bounding boxes.

**Backbone selection:** The SSD in [1] was originally built upon the ImageNet pretrained VGG16 visual classification network [6], whereby the VGG16 was exploited as a high-level feature extractor. Although VGG16 has good object representation capability, it is quite a large network architecture with 24.1 M parameters (cf. Table III). To address this, this work strategically replaces the heavy computing VGG16 with MobileNetv2 [7] (cf. Fig. 1). MobileNet versions of CNNs are lighter architectures due to the usage of depth-wise separable convolution operations. For example, the MobileNet v1 and v2 have nearly a 1/30 of the computational cost and model size as compared to VGG16.

**Multi-scale feature maps:** The standard SSD model uses VGG16 and additional six feature map layers as shown in Fig. 1. They explore features from input images at multiple field of view for detecting objects of various sizes. Table I shows a comparison of the multi-scale object handling with various backbone networks that are considered in this study along with the proposed model. In Table I, we can learn that, the SSD with MobileNet V2 (MB v2) uses feature maps with sizes of  $19 \times 19$ ,  $10 \times 10$ ,  $5 \times 5$ ,  $3 \times 3$ ,  $2 \times 2$  and  $1 \times 1$ , which is different from VGG16-based SSD model. The first feature map layer with the size of  $38 \times 38$  in the VGG16 is shallow and

TABLE I  
THE NUMBER OF FEATURE MAPS USED WITH VARIOUS BACKBONES  
IN THIS WORK FOR ABLATION STUDY.

Feature map layer	SSD with VGG16	SSD with MobileNetV2	Proposed model with MobileNetV2
Layer1	38×38×512	19×19×96	19×19×96
Layer2	19×19×1024	10×10×1280	10×10×1280
Layer3	10×10×512	5×5×512	-
Layer4	5×5×256	3×3×256	-
Layer5	3×3×256	2×2×256	-
Layer6	1×1×256	1×1×256	-

not so effective in extracting key attributes of input images. Therefore, it is modified to have 19×19 feature map.

*B. Proposed Fast Traffic Light & Sign Detection Model*

In SSD with MBv2, it is empirically found that the six feature map layers from 19×19 to 1×1 calculated prior boxes

with a scale of 0.1, 0.2, 0.375, 0.55, 0.725 and 0.9 [1]. Therefore, larger feature maps have prior boxes with smaller scales, so they are ideal for detecting smaller objects. With this institution, we modify the standard SSD. As traffic signs and traffic lights occupy relatively a smaller fraction of the entire road scene, we exploit the first two prediction layers: 19×19, and 10×10 (cf. Table I and Fig. 1), and discard the remaining ineffective layers. This strategic modification saves a lot of memory and computations, resulting in a faster object detection.

*C. Objective Function and Evaluation Metrics*

To train the proposed model, we use a weighted sum of the classification confidence loss ( $L_{conf}$ ) and the localization loss ( $L_{loc}$ ), as given in Eq. (1).

$$L = \frac{1}{N} (L_{conf} + \alpha \cdot L_{loc}), \tag{1}$$

where  $N$  is the number of matched prior boxes, and  $\alpha$  is a constant set to 1. The localization loss is the smooth L1 loss between the predicted and ground truth bounding boxes. Hence, the confidence loss is calculated using a Softmax activation with categorical cross-entropy.

The mean average precision (mAP), which is the primary metric used for comparing the performance of object detectors. The average precision (AP) is computed as the average of maximum precision values at a chosen 11 recall values. Hence, the AP for class  $c$  is defined as in Eq. (2).

$$AP_c = \frac{1}{11} \sum_r \max(P(r)), \tag{2}$$

where  $P(r)$  is the precision for one of the 11 recalls (cf. Eq. 4), and  $r \in \{0.0, \dots, 1.0\}$ . Hence, the mAP for object detection is the average of the APs calculated over all the classes as shown in Eq. (3), where  $C$  is the total number of classes and  $AP_c$  is AP for class  $c$  with IoU (cf. Eq. (5)) of 0.5.

$$mAP = \frac{1}{C} \sum_c AP_c \tag{3}$$

$$Recall = \frac{TruePositive}{TruePositive + FalseNegative} \tag{4}$$

$$IntersectionOverUnion(IoU) = \frac{A_{pred} \cap A_{gt}}{A_{pred} \cup A_{gt}}, \tag{5}$$

where  $A_{pred}$  and  $A_{gt}$  stand for the area of the predicted boxes and ground truth boxes, respectively. An IoU threshold of 0.5 is used to classify whether the prediction is a true positive or a false positive.

IV. EXPERIMENTAL ANALYSIS

*A. Dataset*

This study uses the publicly available road sign detection dataset [25]. As far as we are concerned, currently, there are no peer-reviewed results published on this newly emerging dataset. It consists of 877 images of four classes: speed limit, stop, crosswalk, and traffic light. The dataset is split into training and test sets by taking a ratio of 8:2. Thus, the training and test sets have, respectively, 701 samples with 2103 objects, and 175 samples with 525 objects.

*B. Computational Platform and Training Setup*

The ablation study was carried out on the Google Colab Pro computational platform with the following specifications. A 2.20 GHz Intel(R) Xeon(R) CPU with a 12 GB memory, and a 1.59 GHz NVIDIA T4 GPU with 16 GB memory. The entire program was written in PyTorch 1.9.0+cu102 with Python 3.7.10. The models were retrained using stochastic gradient descent (SGD) optimizer with an initial learning rate of 0.001, momentum of 0.9, weight decay of 0.0005, learning rate decay policy of 10, and batch size is set to 16.

*C. Quantitative Analysis*

To study the impact of proposed framework, four independent investigations, Model A - D, were carried out considering different number of feature map layers with different scales as listed in Table II. Where, Model A uses all the six feature map layers, while Model B was created with only the first 19×19 feature map exclusively omitting other layers. Hence, the Model C was built with only the first two feature map layers, 19×19 and 10×10. The last Model D was built using first three feature map layers: 19×19, 10×10, and 5×5. All the models were fine tuned for 10 epochs. This sanity test shows that the Model C produces the best results and the mAP starts deteriorating when the number of layers is increased to three.

To further validate, these segregated objects' areas are analyzed based on the size of their corresponding bounding boxes with the full image to calculate the percentage of coverage. From Fig. 2, we can see that 95.4% (41.9+53.46) of objects having maximum 20% of image dimension coverage, i.e., majority of the objects come under the range of small/medium size compared to the input image dimensions. That is the reason for increased mAP of the Model C as compared to Model A. Further referring to Table II, one can see that, in terms of the number of trainable parameters, our proposed framework

TABLE II

PERFORMANCE OF THE SSD WITH DIFFERENT NUMBER OF FEATURE MAP LAYER OF THE MOBILE NET V2 (WITHOUT DATA AUGMENTATION)

Model	mAP@IoU0.5	#of Parameters	# of predictions per class
Model A: Full SSD with MBv2	0.307	7,222,518	2,268
Model B: One layer SSD with hMBv2	0.391	3,536,524	1,444
Model C: Two-layer SSD with MBv2	0.433	4,158,146	2,044
Model D: Three-layer SSD with MBv2	0.370	6,932,088	2,194

■ % of objects in the validation set

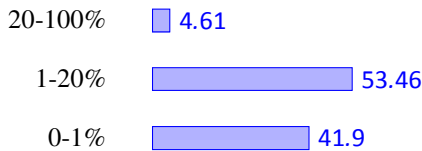


Fig. 2. Percentage coverage of the objects' ground truth bounding box wrt the whole image dimension. The ordinates shows objects' coverage. Based on the % of area cover the bounding boxes they are grouped into 0-1% - small, 1-20%-medium, and 20-100%-large objects.

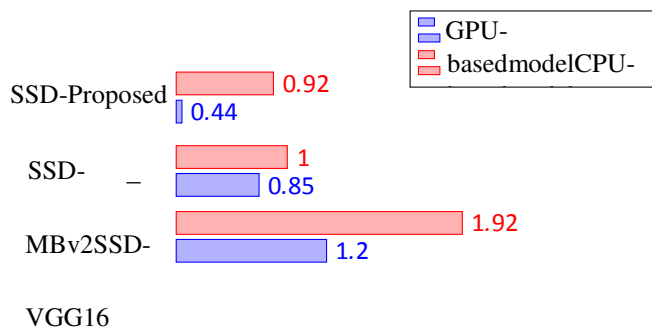


Fig. 3. Detection time analysis of the proposed two-layer SSD with MobileNet V2. The bars show the timing in second.

TABLE III

COMPARATIVE ANALYSIS OF THE PROPOSED MODEL WITH BASELINE SSD MODELS (WITH DATA AUGMENTATION)

Model	mAP@IoU0.5	Detection Time (s)		#of Parameters
		GPU	CPU	
SSD with VGG16	76.7	1.203	1.917	24.1M
SSD with MBv2	68.8	0.848	0.997	7.2M
Proposed SSD with two-layer MBv2	73.8	0.443	0.924	4.1M

with two-layer SSD (Model C) has less number of parameters as compared to Model A. Less number of parameters means a lighter architecture with less computational complexity. Also, the Model C involves 224 lesser number of predictions per class as compared to the Model A with all the six layers leading to a smaller number of computations. All these properties of the proposed solution make it more suitable for real-time detection and deployment of the model on an embedded system.

Using the experimental findings discussed earlier (cf. Table II) as a proof of concept, we finalized the two-layer

SSD with MBv2 and conduct further ablation study. In these experiments, we also apply four data augmentation techniques: photometric distortions, expand image (zoom out), random crop image, and horizontal flip to get better generalization performance. The models are trained with early stopping

(the best mAP is found at 58th epoch). The experimental results are tabulated in Table III, and compared in Fig. 3, and Fig. 4. The average detection times on CPU and GPU are compared for all the three models in Fig. 3. The results show noticeable improvements even in the case of CPU. These results prove that the proposed framework can be deployed on an embedded platform for real-time traffic sign and light detection.

Overall, the proposed two-layer SSD with MobileNet V2 is found to be faster than the baseline SSD with VGG16.

As expected, the SSD with VGG16 has a better mAP due to the deeper feature learning architecture (cf. Table III). However, the proposed model renders a comparable detection performance and achieves a 5% improvement in mAP when compared to the full layer SSD with MBv2. As noticed in Fig. 3, in terms of detection time, the baseline SSD with VGG16 is 63% and 52% slower when compared to the proposed two-layer SSD model, respectively on GPU-based and CPU-based implementations. In a nutshell, the proposed model, nearly 2 times faster than the baseline model with a minor compromise on detection precision < 3%.

#### D. Qualitative Analysis

Fig. 4 shows few visual results for comparing the top-2 models: the proposed two-layer SSD with MB v2 and the baseline SSD with VGG16. We can clearly notice that the proposed model's performance is much better as compared to the baseline model, especially in image IDs: road807, road748, road716 and road213 in terms of small object detection. It is

also noticeable in image ID: road821 that the proposed model can detect two extra traffic lights that the baseline model is not able to detect.

#### V. CONCLUSION

This work presents an efficient exploitation of the SSD model for a fast traffic sign and light detection. In the proposed framework, the standard SSD architecture's backbone, VGG16 is replaced with the lighter MobileNet v2 and only the top-two feature map layers of the SSD are used. Such modifications are proved to be essential not only for small object detection but also quicker detection. The proposed model achieves 2 times (or more) lesser detection time than the baseline SSD models and a comparable detection performance of 76.7% mAP. The proposed solution can be further improved by training on a bigger dataset. Apart from the intended purpose, it can be beneficial for the applications, where the detection of small objects plays an important role.

#### REFERENCES

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

DetectionResultsfromtheStandardSSDwithVGG16



DetectionResultsfromtheProposedModel(SSDwithTwo-layerMobileNetV2)

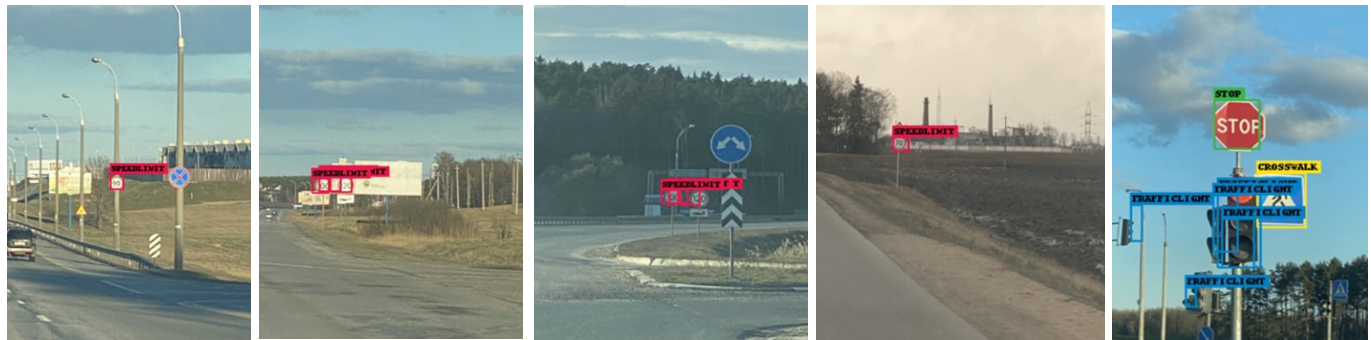


Fig.4.FewSamplesofQualitativeResults.ImageIDsfromCol.1to5:road807,road748,road716,road213,androad821.

- [2] Y. Koo, C. You, and S. Kim, "Opencl-darknet: An openclimplementation for object detection," in *2018 IEEE Intl. Conf. on Big Data and Smart Computing (BigComp)*.IEEE,2018,pp.631–634.
- [3] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: towards real-timeobject detection with region proposal networks," *IEEE trans. on patternanalysisandmachineintelligence*,vol.39,no.6,pp.1137–1149,2016.
- [4] J. Dai, Y. Li, K. He, and J. Sun, "R-fcn: Object detection via region-based fully convolutional networks," in *Advances in neural informationprocessing*,2016,pp.379–387.
- [5] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only lookonce: Unified, real-timeobject detection," in *Proc. of the IEEE Conf. on comput. vis. and patternrecog.*,2016,pp.779–788.
- [6] K. Simonyan and A. Zisserman, "Very deep convolutional networks forlarge-scaleimagerecog." *arXivpreprintarXiv:1409.1556*,2014.
- [7] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proc. of the IEEE Conf. on comput. vis. and patternrecog.*,2018,pp.4510–4520.
- [8] H.-K. Kim, Y.-N. Shin, S.-g. Kuk, J. H. Park, and H.-Y. Jung, "Night-timetrafficlightdetectionbasedonsvmwithgeometricmomentfeatures," *Intl. Journal of comput. and Information Engineering*, vol. 7,no.4,pp.472–475,2013.
- [9] M. Diaz-Cabrera, P. Cerri, and J. Sanchez-Medina, "Suspended trafficlights detection and distance estimation using color features," in *201215th Intl. IEEE Conf. on intell. Transportation sys*.IEEE, 2012, pp.1315–1320.
- [10] M. Omachi and S. Omachi, "Traffic light detection with color and edgeinformation," in *2009 2nd IEEE Intl. Conf. on comput. Science andInformationTechnology*.IEEE,2009,pp.284–287.
- [11] M. Swathi and K. Suresh, "Automatic traffic sign detection and recog.:A review," in *2017 Intl. Conf. on Algorithms, Methodology, Models andApplicationsinEmergingTechnologies*.IEEE,2017,pp.1–6.
- [12] H. Supreeth and C. M. Patil, "An approach towards efficient detectionand recog. of traffic signs in videos using neural networks," in *2016 Intl. Conf. on Wireless Communications, Signal Processing and Networking (WiSPNET)*.IEEE,2016,pp.456–459.
- [13] B. T. Nguyen, S. J. Ryong, and K. J. Kyu, "Fast traffic sign detectionunder challenging conditions," in *2014 Intl. Conf. on Audio, LanguageandImageProcessing*.IEEE,2014,pp.749–752.
- [14] Y. Yang, H. Luo, H. Xu, and F. Wu, "Towards real-time traffic signdetection and classification," *IEEE trans. on intell. transportation sys.*,vol.17,no.7,pp.2022–2031,2015.
- [15] T. Akilan, Q. J. Wu, A. Safaei, J. Huo, and Y. Yang, "A 3d cnn-lstm-based image-to-image foreground segmentation," *IEEE Transactions onIntelligentTransportationSystems*,vol.21,no.3,pp.959–971,2019.
- [16] T. Akilan and Q. J. Wu, "sendec: An improved image to image cnn forforeground localization," *IEEE Transactions on Intelligent TransportationSystems*,vol.21,no.10,pp.4435–4443,2019.
- [17] T. Akilan, Q. J. Wu, A. Safaei, and W. Jiang, "Alatefusionapproachforharnessingmulti-cnnmodelhigh-levelfeatures," in *2017IEEEInternationalConferenceonSystems, Man, and Cybernetics(SMC)*.IEEE,2017,pp.566–571.
- [18] M. Weber, P. Wolf, and J. M. Zöllner, "Deeptlr: A singledeepconvolutional network for detection and classification of traffic lights," in *2016IEEEIntell. vehicle Symposium (IV)*.IEEE,2016,pp.342–348.
- [19] L. B. Karsten, N. Libor, and B. Rami, "A deep learning approach tottraffic lights: Detection," in *2017 IEEE Intl. Conf. on Robotics andAutomation*,2017,pp.1370–1377.
- [20] J. Müllerand K. Dietmayer, "Detectingtrafficlightsbysingleshotdetection," in *2018 21st Intl. Conf. on intell. Transportation sys. (ITSC)*.IEEE,2018,pp.266–273.
- [21] D. Yudin and D. Slavioglo, "Usage of fully convolutional network withclustering for traffic light detection," in *2018 7th Mediterranean Conf. onEmbeddedComputing(MECO)*.IEEE,2018,pp.1–6.
- [22] J. Zhang, M. Huang, X. Jin, and X. Li, "A real-time chinese traffic signdetectionalgorithmbasedonmodifiedyolov2," *Algorithms*,vol.10,no.4, p.127,2017.
- [23] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proc. of the IEEE Conf. on comput. vis. and patternrecog.*, 2017, pp. 7263–7271.
- [24] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu, "Traffic-signdetection and classification in the wild," in *Proc. of the IEEE Conf. oncomput. vis. and patternrecog.*,2016,pp.2110–2118.
- [25] M. Andrew, "Roadsigndetectiondataset," <https://www.kaggle.com/andrewmvd/road-sign-detection>, May2020.