# Vehicle and License Plate Recognition with Novel Dataset for Toll Collection

Muhammad Usama, Hafeez Anwar, Abbas Anwar, Saeed Anwar,

*Abstract*—We propose an automatic framework for toll collection, consisting of three steps: vehicle type recognition, license plate localization, and reading. However, each of the three steps becomes non-trivial due to image variations caused by several factors. The traditional vehicle decorations on the front cause variations among vehicles of the same type. These decorations make license plate localization and recognition difficult due to severe background clutter and partial occlusions. Likewise, on most vehicles, specifically trucks, the position of the license plate is not consistent. Lastly, for license plate reading, the variations are induced by non-uniform font styles, sizes, and partially occluded letters and numbers. Our proposed framework takes advantage of both data availability and performance evaluation of the backbone deep learning architectures. We gather a novel dataset, *Diverse Vehicle and License Plates Dataset (DVLPD)*, consisting of 10k images belonging to six vehicle types. Each image is then manually annotated for vehicle type, license plate, and its characters and digits. For each of the three tasks, we evaluate You Only Look Once (YOLO)v2, YOLOv3, YOLOv4, and FasterRCNN. For real-time implementation on a Raspberry Pi, we evaluate the lighter versions of YOLO named Tiny YOLOv3 and Tiny YOLOv4. The best Mean Average Precision (mAP@0.5) of 98.8% for vehicle type recognition, 98.5% for license plate detection, and 98.3% for license plate reading is achieved by YOLOv4, while its lighter version *i.e.*, Tiny YOLOv4 obtained a mAP of 97.1%, 97.4%, and 93.7% on vehicle type recognition, license plate detection, and license plate reading, respectively. The dataset and the training codes are available at https://github.com/usama-x930/VT-LPR

*Index Terms*—Object detection, Object recognition, Image classification, Machine learning, License plate recognition, Characters recognition, Deep learning

## I. INTRODUCTION

TOLL tax is one of the important means of revenue generation for the departments responsible for the operation and maintenance of highways and motorways. The process of toll tax collection is automated via sensors and cameras in many countries [1] where the license plates and vehicle types are automatically detected to calculate and deduct the toll tax for vehicles. Among others, the main advantage of such automation is to avoid long queues on toll collection stations or toll plazas of highways on occasions such as weekends and national holidays, as shown in Figure 1. However, such systems rely on prior knowledge of the specific format of license plates and their prescribed positions on the front of the vehicles. Due to this reason, their performance is most likely to degrade in places where arbitrary formats and positions of
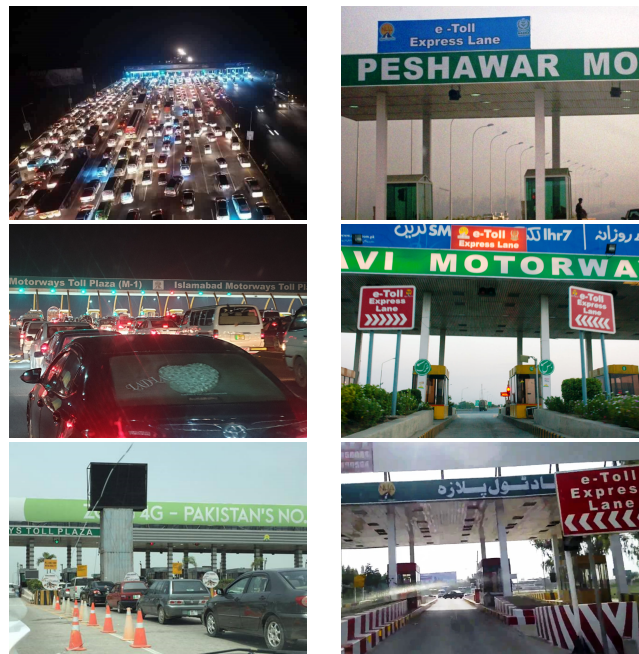
M. Usama, H. Anwar, and A. Anwar are with The Department of Electrical and Computer Engineering, COMSATS University Islamabad, Attock Campus, Pakistan e-mail: hafeez.anwar@cuiatk.edu.pk (Corresponding Author)

S. Anwar is with the Commonwealth Scientific and Industrial Research Organisation (CSIRO), the Australian National University (ANU), and the University of Technology Sydney (UTS), Australia



Fig. 1: Toll plaza traffic congestion[left column] and e-tag based express lane [right column]. These scenarios show the importance of automatic license plate recognition.

license plates are used by most vehicles. An existing solution to this problem is the usage of RFID-based e-tag fitted on vehicles. Figure 1 shows specific lanes on toll plazas where scanners are installed to automatically scan and read the e-tags for tax deduction without stopping the vehicles.

However, this comes with an additional cost of installing extra infrastructure on toll plazas and e-tags on the vehicles. In addition, e-tags are only installed on regular users of motorways that constitute a minor portion of the whole traffic volume.

Leveraging upon the installed surveillance cameras of toll plazas, we propose an image-based solution for toll tax collection. Such a solution is motivated by the recent state-of-the-art results achieved by the usage of convolution neural network (CNN) for image-based vehicle type [2] and license plate recognition [3]. Our framework calculates the toll tax of a vehicle from its image using three steps strategy. The first step deals with the localization and type recognition of the vehicle in the image. The region of the image that depicts the vehicle is then searched for the license plate. The final step is then performed on the localized license plate by reading its characters and digits. All three steps are performed with a CNN-based object localization framework. However, the

| Intra-Class variations | Non-Uniform positions | Multiple license plates | Handwritten license plates | Occluded license plates | Damaged license plates |
|---|---|---|---|---|---|



Fig. 2: Exemplar images of various challenges caused by intra-class type variations due to decorations, license plate positions are not coherent, more than one license plate displayed at vehicle's front, license plates are occluded by decorations and broken or faded license plates.

task of the image-based vehicle and license plate recognition becomes challenging due to variations found in the images of the Pakistani vehicles. Following is the list of challenges whose exemplar images are also shown in Figure 2.

- **Intra-class variations**: Our framework aims at the vehicle type recognition from its front that is affected by the traditional decorations [4] on the vehicles. Such decorations cause dissimilarities among vehicles of the type.
- **Background clutter**: The decorations surrounding the license plate cause background clutter, thus making it difficult to get detected.
- **Non-uniform position**: In the vast majority of the vehicles, the official license plate positioning is not followed, thus making it difficult to detect the license plate.
- **Multiple license plates**: It is also noted that in some vehicles, two license plates are displayed where one is on the front side, and the other is displayed inside the windshield.
- **Non-uniform fonts**: The variations in fonts caused by the use of non-official license plates make it difficult to recognize the characters and digits on the license plates.
- **Partial occlusions**: Due to the dusty environment, the license plates are always partially occluded with dust and mud.
- **Damaged license plates**: In some cases, the non-official license plates are damaged, which causes the missing of some digits or characters.
- **Lack of license plates**: The license plate can even be missing in cases like new vehicles where at the place of

the license plate, a message of "Applied for" is displayed.

Due to these challenges, our proposed framework consists of three individual steps, which are (i) vehicle detection and type recognition, (ii) license plate detection, and (iii) license plate reading. We treat each step as an object detection and recognition problem, where we train a separate model for each of them. Consequently, all three models are applied in a cascaded manner during the run time to recognize both the vehicle and its license plate for calculating the toll tax. In this regard, the following are our novel/significant contributions.

1) We collect a dataset of vehicles that consists of 10K images. This is one of the largest and most diverse image datasets of vehicles to the best of our knowledge.
2) From these images, we generate three sub-datasets each for vehicles type recognition, license plate detection, digits, and character recognition of license plates.
3) We extensively evaluate the performance of several object detection frameworks, which are YOLOv4 and Tiny YOLOv4, YOLOv3 and Tiny YOLOv3, YOLOv2, and FasterRCNN for each of the three steps.
4) For IoT-based and real-time solutions, we also deploy our proposed framework on a Raspberry Pi along with a Pi camera and a graphical user interface (GUI).

The rest of the paper outlines literature review in Section II, dataset in Section III, the proposed methodology in Section IV, experimental results in Section V and conclusion in Section VI.

## II. Literature Review

Since we propose to perform both vehicle and number plate recognition, this section summarizes some literature on both problems.

### A. Vehicle type recognition

Vehicle type recognition deals with the coarse-grained classification of a vehicle where the aim is to put it into one of the classes such as car, bus, van, and truck, *etc.* The image-based vehicle type recognition methods are either based on the so-called handcrafted features or use convolutional neural networks. We give a summary of some of the techniques from both groups.

**Handcrafted features**: The Scale-invariant Features Transform (SIFT) [5] is used by [6] for appearance representation while the relative positions of SIFT feature as a structural feature to deal with the variations found in vehicle images. These features are then utilized via multiple kernel learning for vehicle type recognition. A discriminative shape descriptor is developed with the help of a modified SIFT and edge points for vehicle type recognition [7]. In a similar manner, vehicle type recognition is performed using SIFT by other methods [8], [9], [10]. Sun *et al.* [11] employ a two-step classification strategy where in the first step, the detected vehicle is classified as a small or large vehicle by using global features generated with improved Canny Edge detection. In the second step, the finer type is predicted for the vehicle via Gabor wavelet kernels extracted at five scales with eight orientations.

Peng *et al.* [12] utilize the location of the license plate as a prior to extracting the front of the vehicle. The feature vector of the extracted vehicle front is then composed with the help of Eigenvectors. K-means clustering is then applied over these vectors in order to represent a vehicle type with a cluster center. A given test sample is then classified by calculating the distance of its feature vector to that of a cluster center. A follow-up work [13] improves this method by enriching it with other information queues such as type-specific license plate color and background subtraction. Others handcrafted features-based for vehicle type recognition use gradient and edges [14], edge direction-based deformable templates [15], and Gaussian Mixture Model [16].

**CNN-based methods**: A two-steps strategy is proposed [17] for vehicle type recognition, where in the first step, class proposal regions are generated for vehicle localization. In the second step, descriptors are obtained for these regions using CNN embeddings, where the vehicle type recognition is then performed using an SVM over these embeddings. The vehicles types used in the study include cars, SUVs, single and double trailer trucks, where for cars and SUVs, the achieved precision is 95%, and for the trucks, it ranges between 92% and 94%. Vehicle type recognition is performed [18] in surveillance videos using CNN. However, the training image data comes from the web, due to which the authors modify the objective function to get a generalized transfer learning for vehicle type recognition. In order to reduce the cost of manual annotation of the surveillance videos, A performance evaluation of the traditional method of SVM on SIFT features against a deep neural network (DNN) is done by [19] DNN performs the best.

Kim & Lim [20] train various CNN models on randomly sampled images from the training set. The number of sampled images is equal to half the size of the training set. The so-called active learning technique is used by [21] where the network training is started with relatively smaller and labeled training data. Further data is then automatically selected by the network, both with high and low entropy for the retraining process. However, in training, data augmentation is also used to generate the augmented images by flipping and rotating the original photos. Each trained model generates a prediction for a test image, where a weighted voting process calculates the final prediction. A similar strategy [22] is adopted, where the ensemble of the so-called local and global networks is used such that the local networks are trained on subsets of training set while the global is trained on the entire training dataset. The inference for a test image where a single prediction is selected from local networks. The final outcome is then made by aggregating the predictions of all global networks and the ones chosen from the local networks.

Since, at toll tax collection, the vehicles line up properly at the toll plazas, we aim to recognize a vehicle from the front. Nonetheless, our approach to performing image-based vehicle type recognition is different from all the previously stated approaches in the following manner.

- We face the problem of severe clutter caused by the traditional decoration measures.
- Such decorative measures also cause high intra-class variations, such that vehicles belonging to the same class look very different from one another.
- Finally, the high dusty conditions at highways in Pakistan also cause a major challenge for vehicle type recognition

### B. License plate detection and recognition

The second important step of our proposed pipeline is the localization of the license plate and then recognizing the digits and characters. From a detected vehicle in the image, the aim is to detect, localize, extract and then read its license plate. Vehicle license plate detection and reading is a well-researched problem; however, it is still an active area of research. We give an overview of both the handcrafted features-based and CNN-based methods in the following.

**Handcrafted features**: Farmanullah *et al.* [23] perform detection of the license plates using image processing operators. The character recognition on the detected license plates is then performed by evaluating various machine learning classifiers over the geometric features of the characters. The Histogram of Oriented Gradients (HoG) [24] with a Support Vector Machine (SVM) in order to detect and recognize license plates in challenging weather conditions and uneven illumination [25]. Other methods of license plate recognition use SIFT [26] [27], Maximally Stable Extremal Regions (MSER) [28] [29] and a combination of MSER and SIFT [30].

**CNN-based methods**: The CNN-based framework [31] is proposed in which, before localizing LP, a vehicle is detected to eliminate false positives and other objects such as signboards, *etc.*, employing CNN-based famous object detector

| Buses | Cars | Carry vans | Truck Type 1 | Truck Type 2 | Vans |

Fig. 3: Exemplar images of vehicle types. The intra-class variations can be observed due to decorations, manufacturers and other shape modifications such bumpers and steel bodies for goods transport

framework YOLO [32] to localize the LP. Various versions of Bangladeshi license plates [33] are detected using a modified version of YOLO called YOLOv4, employing Tesseract [34] as an OCR engine to recognize characters from the detected license plates. Mask RCNN [35] is used for the detection, segmentation, and recognition of Tunisian license plates [36] under challenging image variations caused by environmental conditions, cluttered background, orientation, and language differences.

Yet another CNN-based method [37] is proposed for the detection and recognition of license plates. However, unlike most techniques, the authors relied on synthetically generated license plate images that suffered from the common variations found in the real license plate images. Consequently, showing that their framework outperformed the CNN trained on natural images of the license plates. Yang *et al.* [38] propose a real-time coarse to fine strategy via contours detection and reconstruction for license plate detection while using a CNN for final character recognition. A layout-independent method [39] is proposed for detecting and recognizing the license plates that belong to several countries. To this end, as a first step, detecting vehicle(s) in a given image and then detect and read the license plates. Both the steps are performed using the YOLO V3 object localization framework.

Wang *et al.* [3] propose a two-step process for license plate detection and recognition by designing two specialized networks. For detection, VertexNet is designed to extract the spatial information of the license plate and later on uses this information to rectify the detected license plate image. For character recognition, the detected license plate is then fed to the second network named SCR-Net. The semantic

segmentation based on DeepLabv2 [40] is utilized by [41] to extract the license plate from the image, which is then followed by counting the number of segments that are candidates of license plate characters. Such segmentation and counting make their method fast and accurate as of the individual steps of license plate detection, extraction, and recognition are bypassed. Other CNN-based methods for license plate detection and recognition use YOLO [42], LSTM [43], and Attentional Networks [44]. Nonetheless, we aim at license plates detection and recognition such that, unlike most previous methods, their locations on the vehicles, layout, and font styles are not standardized.

## III. DATASET

We collect a novel dataset, *Diverse Vehicle and License Plates Dataset (DVLPD)*, that has 10k images of Pakistani vehicles. These vehicles are trucks, buses, vans, carry vans, and cars. There are two sub-classes of trucks called *type1* trucks having a single axle and *type2* trucks that have double axles. Figure 3 shows the exemplar images for each vehicle. Since our approach to the problem is data-driven, we carefully accommodate all the challenges mentioned above in the images of each vehicle type. After dataset collection, the pre-processing steps and ground truth generation are performed to train and test the CNN models. These steps are explained subsequently and are shown in Figure 4, whereas the statistics about vehicles images and their license plates are given in Figure 5.

Fig. 4: Ground truth generation process the three steps.

## A. Data pre-processing
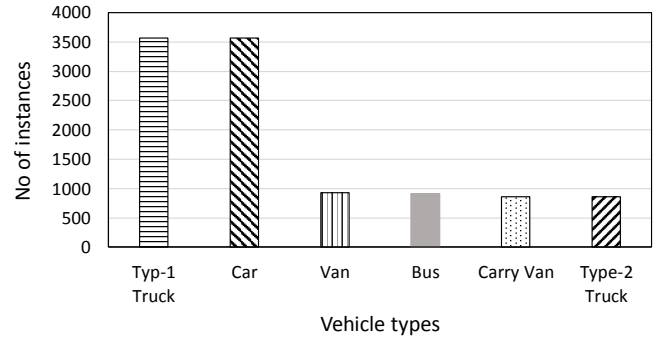
1) **Cropping**: The images are cropped with a standard aspect ratio of 3:4 to remove the unwanted information. For instance, a vehicle image taken from a long distance contains a patch of road that acts as redundant information. Apart from this, cropping also reduces the image size, which is important for the reduced training time of the deep learning model.

2) **Resizing**: The images are resized because the network must receive the same size images during the training process. The image size depends greatly on the problem being dealt with and the architecture's requirement. Hence, there is no standard defined for the size of the input images; it can be of different resolutions, such as $416 \times 416$ and $256 \times 256$. Images with lower resolution require less processing time, but the trade-off is fewer details in the image. In this research, the image is resized to $416 \times 416$ pixels for the training and testing. Most architectures recommend the chosen size because the details are preserved without losing important information in the images.

## B. Ground truth generation

As mentioned previously, we train a separate model for each of the three steps *i.e.*, vehicle type recognition, license plate detection, and license plate character recognition. Consequently, all the vehicle images are annotated for each of the individual steps in the following manner.

1) **Vehicle type detection**: First, we annotate all the pre-processed images for six vehicle types using the LabelIMG tool [45]. The ground truth bounding box is manually drawn only around the depicted vehicle in the image. This labeled dataset is employed to train and test the model for vehicle type recognition.

2) **License plate detection**: Subsequently, the ground truth bounding boxes around the vehicles are used to crop out all the depicted vehicles. This generates a sub-dataset which we call "*lp-detect dataset*". The dataset images are then prepared for the second step of our framework, which is license plate detection. The bounding boxes are manually drawn around the license plates of the vehicles. These annotated images are utilized for training and testing the second model for license plate detection.

3) **License plate reading**: For the final step of license plates reading, we make another dataset by cropping out

(a) Number of instances for vehicle types in *pre-processed* and *lp-detect dataset*



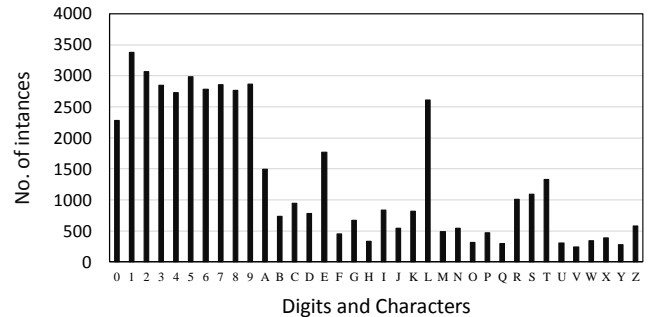(b) Number of instances for each digit and character in *lp-read dataset*



Fig. 5: Statistics of the *pre-processed*, *lp-detect dataset*, and *lp-read dataset dataset*

all the license plates from the *lp-detect dataset dataset* which is named as *lp-read dataset dataset*. This dataset is then annotated for digits (0-9) and characters (A-Z) for the final step of license plates reading.

The statistics of each of the datasets are shown in Figure 5.

## IV. METHODOLOGY

The proposed framework for the image-based toll tax calculation is explained in this section. The toll tax is calculated by designing a framework that performs image-based end-to-end recognition of the vehicle's type and license plate. Due to the challenging image variations, we have focused on developing a system that uses three object detection models. With such
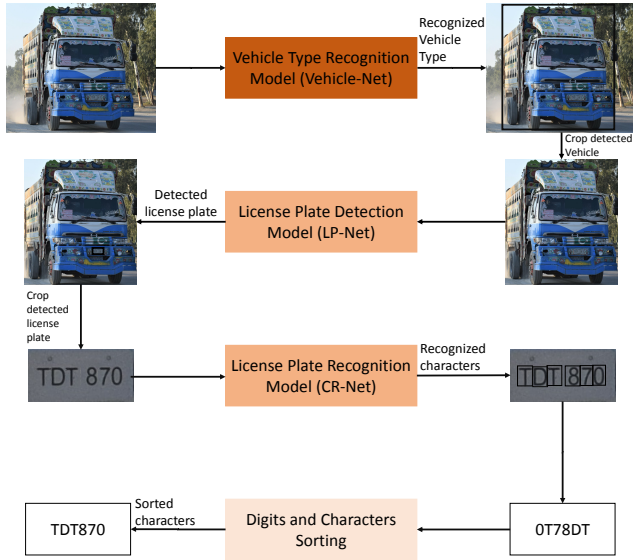
Fig. 6: Block Diagram showing the full working for our proposed framework. The different colors represent a specific module in our framework.

a strategy, we aim to achieve the highest possible accuracy at recognizing the type of the vehicle, localize its license plate, and finally identify the license plate's characters. Each of these steps is explained in the following and their deployment on a Raspberry Pi for real-time usage.

### A. Proposed Model

Figure 6 shows the block diagram of the entire framework. The input to this framework can be an image, series of images, or video frames, whereas the outputs are the recognized vehicle type and license plate characters. When an image is fed to the framework, the first object detection model is triggered for the vehicle type recognition, which we call the "*Vehicle-Net.*" The vehicle is localized in the image via a bounding box, and its type is predicted. The coordinates of this bounding box are then used to extract the vehicle, and the next model for license plate detection called the "*LP-Net*" is applied to it. This model predicts the bounding box around the license plate of the vehicle. Finally, the last model called the "CR-Net" is then triggered on the detected license plate, which specializes in detecting and recognizing the characters of the license plate. In such a way, each input frame is processed seamlessly by the three models, with a final output giving the labels of the detected type of vehicle and the characters in its license plate. When there is no vehicle present in the input frame at the first object detection network, the frame is skipped, and the system waits for the next frame to detect a vehicle in the image.

### B. RaspberryPi Deployment

After the offline training, all three models are deployed on a RaspberryPi 4 (RPi). As shown in Figure 5, after acquiring images with the help of a Pi camera interfaced at the camera port, the device applies all three models on them for a
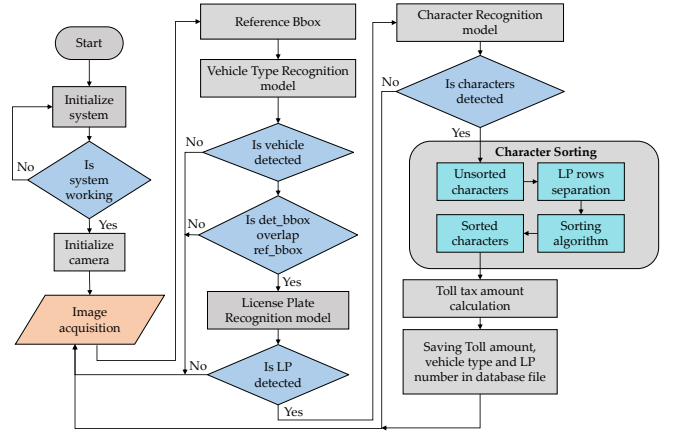


Fig. 7: Flow Process of Proposed Methodology

vehicle type and license plate recognition. Figure 7 shows the flowchart that depicts the technical working of the algorithm running on the RPi. As a first step, the models are loaded into the memory of RPi. The RPi has a RAM of 2GB which is sufficient to load the tiny models of YOLOv3 and YOLOv4. In the next step, the camera is initialized, and a video stream is an input into the models' frame by frame. In the next step, a reference bounding box is drawn to tackle the problem of multiple detections, which is explained in Section IV-B1. After the vehicle is detected, the output of the first network is cropped out using the coordinates of the bounding box with the highest confidence score. This output is then fed into the next model to detect the license plate in the frame. Similarly, the frame is processed again from the network's final layer, and this time the license plate is fed into the last network, which recognizes the characters in the image. The recognized characters are then sorted according to their detected bounding boxes using the sorting algorithm, which is explained in Section IV-B2.

*1) Multiple Vehicle Detection:* Multiple vehicle detection in a single image makes it difficult to select the vehicle for which toll tax calculation must be performed. The camera position on the toll tax collection plaza gives a rough estimate of the expected position of the vehicle in the image region. We draw a so-called "*reference box*" around this image region and store its coordinates. Once a vehicle is detected, the algorithm iterates over all the bounding boxes detected by the network and checks whether the overlap area between the reference and detected bounding box exists using a formula on the coordinates of both reference and detection. The detailed algorithm for avoiding multiple vehicle detections is outlined in Algorithm 0, where the coordinates of the reference bounding box are compared with the detection results.

*2) Character Sorting in License Plate:* The sequence of the recognized characters and digits is not identical to that of the license plate due to the fact that *CR-Net* does not sort its output. While dealing with license plates, the sequence of the characters is vital. The characters of license plates are usually arranged in one or two rows, as shown in Figure 8. To sort the characters in more than one row, we use the detected labels and the coordinates of their bounding boxes to sequence the

---

**Algorithm 1** Multiple Vehicle Detection

---

**Require:** Draw Reference Bounding Box on camera frame
**Ensure:** $InputFrame \geq 0$
  $RefBBox[\ ] \leftarrow$ *xmin, ymin, xmax, ymax*
  $indices \leftarrow$ *Detection Results*
  **for** i in indices **do**
    **if** $len(indices)$ is $> 0$ **then**
        $left \leftarrow$ *Detected BBox Left Coordinate*
        $right \leftarrow$ *Detected BBox Right Coordinate*
        $top \leftarrow$ *Detected BBox Top Coordinate*
        $bottom \leftarrow$ *Detected BBox Bottom Coordinate*
        $DetBBox[\ ] \leftarrow$ *xmin, ymin, xmax, ymax*
        $xA \leftarrow (max(RefBBox[0], DetBBox[0]))$
  ▷ Computing max and min of RefBBox[ ] and DetBBox[ ]
        $yA \leftarrow (max(RefBBox[1], DetBBox[1]))$
        $xB \leftarrow (max(RefBBox[2], DetBBox[2]))$
        $yB \leftarrow (max(RefBBox[3], DetBBox[3]))$
        $overlapArea \leftarrow abs(max((xB - xA, 0)) *$
$max((yB - yA), 0))$
        **if** $overlapArea \leq 0$ **then**
            Overlap Exists
        **else**
            No Overlap Exists
        **end if**
    **end if**
  **end for**

---

characters in the license plate correctly. The labels and the coordinates of those labels are stored, then the bounding boxes with labels as alphabets and numbers are separated. Then, the coordinates of the bounding boxes in the alphabets and numbers are sorted using the "x-min" coordinate from left to right in detected bounding boxes and displayed as the final correctly sequenced license plate.

*C. Backbone Architectures*

Various State-of-the-art (SOTA) architectures are used as the backbone model in our proposed methodology to accomplish vehicle type recognition, license plate extraction, and character recognition tasks.

1) **Faster RCNN**: Faster RCNN uses Deep CNN called Region proposal Network (RPN) to generate the sets of region proposal bounding boxes. After the region proposals are obtained, they are subjected to another network for extracting the proposal features from the feature map and classification using softmax. This belongs to the RCNN family, which is a multi-stage object detection algorithm due to which the real-time performance is barely reachable.

2) **YOLOv2**: Single staged end-to-end training and detection are achieved through YOLO. This one-stage detector treats object detection as a regression problem, thus closing the gap of real-time performance. YOLOv2 framework is called Darknet-19 which consists of 19 convolutional layers and five max-pooling layers. It uses multi-resolution images by removing fully connected layers. Global average pooling is used for the prediction, and batch normalization is introduced after each convolution layer to regularize the model and speed up the convergence.

3) **YOLOv3 & Tiny YOLOv3**: YOLOv3 is more optimized in terms of detection speed and accuracy as compared to previous models. Feature Pyramid Network (FPN) approach is introduced in YOLOv3, which allows detecting the objects at three different scales. It also uses residual blocks to tackle the problem of vanishing gradients in a dense network. The YOLOv3 tiny is a relatively lightweight architecture as compared to the YOLOv3. It is utilized for an embedded system where the resources are limited. The detection accuracy of the lighter architecture is less than the complete architecture, while the performance in FPS for Tiny YOLOv3 is greater than YOLOv3. The feature extractor is the main difference in both architectures, where YOLOv3 uses the Darknet-53 with 53 layers, and the Tiny YOLOv3 uses a compressed version with seven convolutional layers and six max-pooling layers in the feature extractor.

4) **YOLOv4 and Tiny YOLOv4**: YOLO v4 is introduced with numerous improvements, in particular with Path Aggregation Network (PAN). It is used to communicate information from lower layers to higher ones. In YOLO v4, the CSPDARKNET53 model is used as the backbone, a CNN that uses DarkNet-53. It employs a split and merges strategy to split the feature map of the base layer into two parts and then merge them. This allows for more gradient flow through the network. Spatial Pyramid Pooling (SPP) and PAN are employed as the neck and YOLO v3 network as the head. A new data augmentation technique is applied called Mosaic data augmentation. In this technique, four images of different aspect ratios are combined into a single image, which helps the model to find smaller objects and pay less focus to the surrounding scenes. Tiny YOLOv4 is a lighter and compressed version of YOLOv4. Tiny YOLOv4 has a simpler architecture that contains 29 convolutional layers, whereas YOLOv4 contains 137 pre-trained convolutional layers; thus, the tiny version is more feasible for embedded and mobile applications. Tiny YOLOv4 is eight times faster in terms of Frame Per Second (FPS) than YOLOv4, whereas accuracy for tiny version is 2/3rds of YOLOv4.

## V. RESULT AND DISCUSSION

The experimental results are discussed in this section systematically. In the first place, we outline the protocols and hardware used in the experiments. We then give some attributes of the datasets used in the training and testing of each model. The results for vehicle type recognition, license plate detection, and license plate recognition are then reported in the form of precision, recall, F1-score, and mean average precision (mAP). We then report an end-to-end accuracy achieved by our framework, and finally, the execution time of the backbone architectures on Raspberry Pi is reported.

Fig. 8: Different types of license plates showing the diversity in style and arrangement.

TABLE I: Hardware Specifications for training the detection models

| Resource Type | Specification |
|---|---|
| CPU | Intel Core i7-6700 |
| RAM | 16GB DDR 4 |
| GPU | Nvidia GeForce GTX 1060 |
| CPU Cores | 4 |
| Frequency | 3.4 GHz |

*A. Experimental Protocols*

For the training and testing of each model, we randomly split their related image dataset into 80% training and 20% test set and recorded their results. Since the split is random, we carry out such training and testing five times and then report their average results and standard deviation. The hardware system specifications for training and testing are shown in Table I.

We use the Darknet framework to implement YOLOv4, Tiny YOLOv4, YOLOv3, Tiny YOLOv3, YOLOv2, and the TensorFlow Object Detection API for faster RCNN implementation. The trained models are tested with the color images of resolution $416 \times 416$ pixels. The YOLO model has three scale pyramids used to detect small, medium, and large objects. The number of epochs for training Vehicle-Net, LP-Net, and CR-Net is 8k, 2k, and 10K, respectively. The batch size is 64, where the subdivisions are 32. The activation function used in YOLOv3, Tiny YOLOv3, and YOLO v2 is Leaky ReLU, YOLOv4, and Tiny YOLOv4 use the Mish activation while FasterRCNN uses the Relu activation function. The training loss depends on the losses caused by objectness, classification, and coordinate loss. The learning rate is fixed at $10^{-3}$. Augmentation is also applied to the training images to accommodate for the variations caused due to saturation, hue, and exposure. The threshold for non-maxima suppression (NMS) is set to 0.3 while the IOU threshold is set to 0.5, which means that the detection is considered correct if the overlap is greater than the confidence score of 0.5. All the parameters are summarized in Table II.

TABLE II: Training parameters of detection models

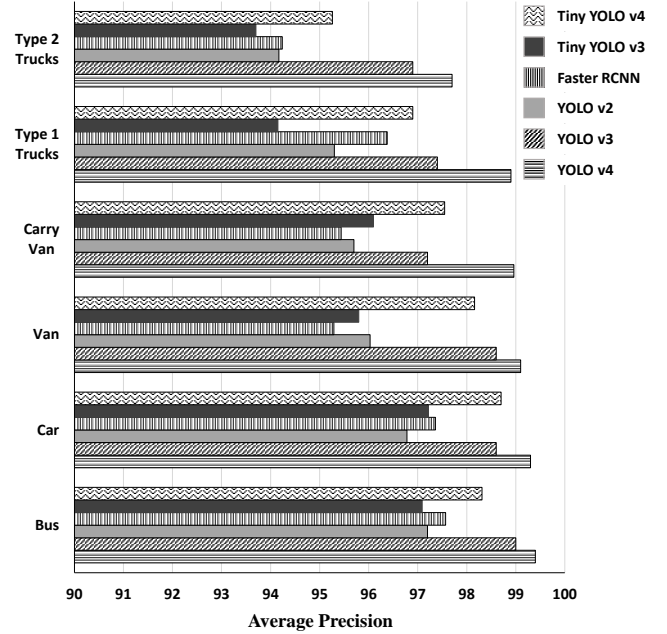| Parameters | Configuration |
|---|---|
| Framework | Darknet & Tensorflow object detection API |
| Max Batches | 3k, 5k & 10k |
| Batch | 64 |
| Subdivisions | 32 |
| Width | 416 |
| Height | 416 |
| Channels | 3 (RGB) |
| Policy | Steps |
| Learning rate | 0.001 |
| Momentum | 0.9 |
| Decay | 0.0005 |
| Saturation | 1.5 |
| Exposure | 1.5 |
| Hue | 0.1 |
| NMS Threshold | 0.3 |



Fig. 9: Average precision achieved for each of the vehicle types

*B. Dataset Attributes*

Our three image datasets are *pre-processed dataset* for Vehicle-Net, *lp-detect dataset* for LP-Net and *lp-read* dataset for the CR-Net. The *pre-processed dataset* training and testing sets contain 7,788 and 1,947 images, respectively, while the total vehicle annotation in this dataset is 10,218. In the *lp-detect dataset* the training and test sets contain 7,891 and 1,972 images, respectively, while the dataset has 9,935 annotated license plates of each vehicle. Lastly, the *lp-read* dataset, which is used to recognize license plate characters, contains 6,261 images in the training set and 1,565 images in the test set. The total number of annotated characters in the dataset is 47,711.

TABLE III: Results of vehicle type recognition achieved by various architectures

| Vehicle Type Recognition | | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | mAP@0.5 |
| **Tiny YOLOv4** | 95.8±1.2 | 96.4±0.5 | 95.8±0.7 | 97.1±0.4 |
| **Tiny YOLOv3** | 92.0±2.3 | 92.2±1.7 | 92.0±1.4 | 94.8±1.0 |
| **Faster-RCNN** | 96.0±0.6 | 95.8±1.9 | 95.6±1.0 | 96.1±0.6 |
| **YOLOv2** | 92.4±1.3 | 92.2±2.3 | 93.2±0.7 | 95.0±0.6 |
| **YOLOv3** | 95.6±1.8 | 97.6±0.5 | 96.4±1.0 | 97.6±0.7 |
| **YOLOv4** | 98.4±0.8 | 97.8±0.4 | 98.2±0.7 | 98.8±0.4 |

TABLE IV: Quantitative performance of state-of-the-art methods for license plate detection.

| License Plate Detection | | | | |
|---|---|---|---|---|
| | Precision | Recall | F1-score | mAP@0.5 |
| **Tiny YOLOv4** | 96.6±1.0 | 97.6±1.0 | 96.8±0.7 | 97.0±0.5 |
| **Tiny YOLOv3** | 93.0±1.4 | 92.0±1.4 | 92.6±1.0 | 95.0±0.7 |
| **Faster-RCNN** | 96.2±1.3 | 96.4±1.0 | 96.2±0.7 | 96.0±0.7 |
| **YOLOv2** | 93.4±1.3 | 93.0±1.7 | 93.4±1.0 | 95.3±0.8 |
| **YOLOv3** | 97.6±0.5 | 97.6±0.5 | 97.6±0.5 | 98.0±0.4 |
| **YOLOv4** | 97.4±1.0 | 97.4±1.2 | 97.4±0.5 | 98.5±0.3 |



Fig. 10: Vehicle type recognition. The last row shows the detection results of the images taken at the toll plaza

### C. Vehicle Type Recognition

The average precision achieved by each architecture for each of the vehicle types is shown in Figure 9. YOLOv4 achieves the highest scores for all the types, while its tiny version performs better than the tiny version of YOLOv3. As mentioned before, due to random shuffling of the dataset, the experiments are performed five times where for each experiment, the number of training iterations is 8K. It can be noted that the best result of 99% is achieved for buses, cars, and vans, while trucks are recognized the least. We believe that this is mainly due to the high intra-class variations in trucks due to the heavy decorations. The least precision is achieved for type-2 trucks, mainly due to their similarity with type-1 trucks due to front decorations and models from similar manufacturers. Nonetheless, YOLOv4 achieves the best mAP of 98.4% for vehicle type recognition followed by YOLOv3 and Tiny YOLOv4 as shown in Table III. Some vehicle detection results are shown in Figure 10 where the last row

shows the detected vehicle on a toll collection plaza.

### D. License plate detection

License plate detection is the most crucial step in our proposed framework. This is since the Pakistani vehicles are heavily decorated and exposed to all sorts of weather conditions. Over time, license plates are covered with dust and mud, which makes their detection challenging. Like the vehicle type recognition, we perform the experiments five times, whereby for each experiment, the *lp-detect* is randomly split into training and test sets. Each time, the number of training iterations are 2K. Figure 11 shows the average precision achieved by each architecture on the license plates of all the vehicle types, while the detailed results are shown in Table IV. The license plates of buses are most accurately detected, while those on the type-1 trucks are detected the least. This is due to several factors; for instance, the positions of the license plates on buses are more coherent than those on trucks. Secondly, the decorations again play a degrading role in license plate detections as trucks; specifically, type-1 trucks are more decorated than buses, cars, and vans. Third, most
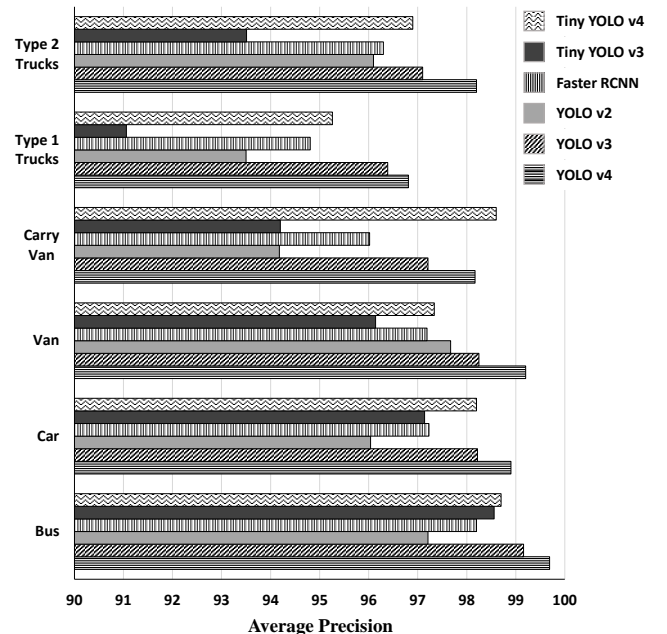


Fig. 11: Average precision achieved by various architectures for the detection of license plates at each of the vehicle types

Fig. 12: License Plate Detection result. The last row shows the detection results on the images taken at the toll plaza

Fig. 13: YOLOv4 average precision ($x$-axis) for each of the digits and characters classes ($y$-axis)

of the type-1 trucks do not display the official license plates, due to which their fonts and styles are not coherent, hence resulting in a lower detection rate. Some license plate detection results achieved by YOLOv4 are shown in Figure 12, where the last row shows the detections on a toll plaza. It can detect license plates with extreme background clutter, such as those of the heavily decorated trucks. Interestingly, damaged license plates and those displayed behind the vehicle's windshield are also accurately detected. Furthermore, it is also able to detect license plates in dusty conditions. The best overall mAP of 98.5% is achieved by YOLOv4 followed by YOLOv3 and Tiny YOLOv4.

### E. License plate characters and digits recognition

The recognition of characters and digits on Pakistani license plates is a challenging task due to the variations found in their fonts and the license plates layouts. As they are depicted on the license plates, their depictions are suffered from dust and mud, thus making their image-based recognition non-trivial. We train and test the third model named "*CR-Net*" 5 times, where each time the dataset is randomly split. For each of the experiments, the number of training iterations is 10K due to
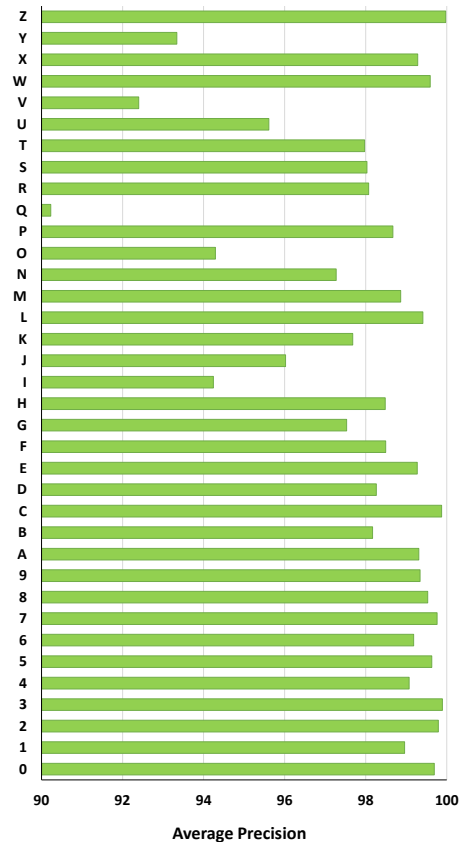
TABLE V: Character recognition results achieved by various architectures

| | Character Recognition | | | |
|---|---|---|---|---|
| | **Precision** | **Recall** | **F1-score** | **mAP@0.5** |
| **Tiny YOLOv4** | 94.0±2.0 | 94.4±2.2 | 94.4±1.5 | 93.8±0.5 |
| **Tiny YOLOv3** | 89.8±1.2 | 92.0±1.7 | 91.4±1.0 | 92.4±1.0 |
| **Faster-RCNN** | 93.0±1.4 | 93.2±1.6 | 93.2±1.5 | 93.0±1.0 |
| **YOLOv2** | 91.6±1.3 | 91.8±1.7 | 91.8±1.1 | 92.8±0.7 |
| **YOLOv3** | 97.6±0.8 | 95.4±1.6 | 96.8±1.0 | 97.0±0.4 |
| **YOLOv4** | 97.2±0.7 | 97.0±1.1 | 97.0±0.6 | 98.3±0.5 |

the relatively higher number of classes in the *lp-read dataset* which is 36. The results are given in Table V, whereas the average precision achieved by the architectures for each of the character's classes is shown in Figure 13. YOLOv4 performs the highest mAP of 98.3%, followed by YOLOv3 and then Tiny YOLOv4. Figure 14 shows some license plate character and digits prediction outputs. As can be seen, our model robustly and effectively tackles the different inconsistencies and variations present in the character dataset. Blurriness, angle shot, different font and font sizes, and resolution do not affect the performance and accuracy of our model to a large extent.

Fig. 14: Character and digits recognition results on the license plate. The output of the framework is shown on the top of each image.

TABLE VI: End-to-end recognition rates for each vehicle type

| Vehicle type | Recognition Rates |
|---|---|
| Car | 99% |
| Van | 98% |
| Carry/Van | 98% |
| Bus | 97% |
| Type2-Truck | 97% |
| Type1-Truck | 96% |

*F. End-to-End recognition*

Nonetheless, we aim for an end-end system for toll tax collection that can perform vehicle type and license plate recognition with the highest possible accuracy. For this purpose, we collect another test dataset where the images are taken at a toll tax plaza. This dataset contains 100 images per vehicle type, and we call it "*end-to-end dataset*". We provide each of the images to the sequential pipeline formed by *Vehicle-Net*, *LP-Net* and *CR-Net*. The backbone architecture used for each model is YOLOv4 due to its best performance on individual tasks. The end-to-end recognition is considered correct only if the vehicle type and license plate are correctly recognized. The results achieved for each of the vehicle types for an end-to-end recognition are shown in Table VI. As expected, cars and vans have the highest accuracy because they are not heavily decorated compared to trucks. Secondly, the positions of license plates on these vehicles are coherent. Lastly, on most of the cars and vans, the official computerized license plates are used, due to which there are fewer variations in their fonts and designs. As expected, both types of trucks have the least end-to-end recognition rates due to their heavy decorations and variations in license plates positions, fonts, and designs.

TABLE VII: The computational time of YOLOv4 and Tiny YOLOv4 on GPU along with Raspberry Pi for all the three models

| | Execution time (sec) | | | |
|---|---|---|---|---|
| | GPU | | Raspberry Pi | |
| | YOLOv4 | Tiny YOLOv4 | YOLOv4 | Tiny YOLOv4 |
| **Vehicle-Net** | 1.2 | 0.16 | 7.1 | 0.95 |
| **LP-Net** | 0.9 | 0.14 | 6.5 | 0.8 |
| **CR-Net** | 0.8 | 0.13 | 5.0 | 0.7 |

*G. Execution Time*

We aim to deploy the proposed framework on a Raspberry Pi to collect the toll tax on plazas. Due to this reason, the images acquired by the Pi camera should be efficiently processed for a vehicle type and license plate recognition in order to calculate the toll tax in real-time. In this regard, we perform experiments to estimate the inference time taken by our proposed framework with both YOLOv4 and Tiny YOLOv4 as its backbone architectures. The comparative time analysis is carried out on Nvidia GTX1060 GPU and Raspberry Pi, and results are shown in Table VII. As expected, due to the lighter architecture, Tiny YOLOv4 takes less time on both GPU and Raspberry Pi than YOLOv4.

*H. Graphical user interface for real-time deployment*

Figure 15 shows the user-friendly graphical interface of our proposed toll tax collection framework. We have given the flexibility of using both images and videos. The threshold is set to 0.5 by default that can be changed depending on the environment parameters. Lastly, outputs of all three steps *i.e.*, vehicle recognition, license plate detection, and character recognition are shown pictorially and in text.

## VI. CONCLUSION

We present a framework for image-based automatic toll tax calculation of Pakistani vehicles, including cars, buses, carry vans, vans, and trucks. This framework can be instrumental in avoiding traffic jams on toll tax collection plazas and can prove an economical alternative to RFID-based systems. Our proposed framework is based on a three-step strategy involving vehicle detection and recognition in images, license plate detection on the detected vehicle, and finally, license plate digits and characters recognition. We evaluated variants of YOLO (v2, v3, and v4), their lighter versions (Tiny YOLOv3 and Tiny YOLOv4), and FasterRCNN for all the three steps on an image dataset of 10K images of Pakistani vehicles on roads and in toll plazas. Overall, YOLOv4 performed the best on all three steps individually and then combined on an end-to-end recognition that spans all three models sequentially. Its lighter version, *i.e.*, Tiny YOLOv4, performed second best; however, it took less time than YOLOv4 on GPU and Raspberry Pi.

Possible future directions of the current research are more focused on trucks. We get the least average precision on them due to the variations in their shapes and models. We also plan

Fig. 15: Graphical User Interface (GUI) for real time deployment of our proposed framework.

to evaluate more object detection frameworks to get better mAP on real-time images and videos.

## REFERENCES

[1] K. TrafficCom, "Tolling," https://www.kapsch.net/en/ktc/solutions/tolling, last accessed: July 30, 2021.

[2] G. K. Chaitanya and G. Maragatham, "Object and obstacle detection for self-driving cars using googlenet and deep learning," in *Artificial Intelligence Techniques for Advanced Computing Applications*. Springer, 2021, pp. 315–322.

[3] Y. Wang, Z.-P. Bian, Y. Zhou, and L.-P. Chau, "Rethinking and designing a high-performing automatic license plate recognition approach," *IEEE Transactions on Intelligent Transportation Systems*, 2021.

[4] S. Ali, "Pakistani truck art," https://gulkhan.pk/, last accessed: July 31, 2021.

[5] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.

[6] Z. Dong and Y. Jia, "Vehicle type classification using distributions of structural and appearance-based features," in *2013 IEEE International Conference on Image Processing*, 2013, pp. 4321–4324.

[7] X. Ma and W. E. L. Grimson, "Edge-based rich representation for vehicle classification," in *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, vol. 2. IEEE, 2005, pp. 1185–1192.

[8] L. Dlagnekov and S. J. Belongie, *Recognizing cars*. Department of Computer Science and Engineering, University of California ..., 2005.

[9] M. Conos, "Recognition of vehicle make from a frontal view," Ph.D. dissertation, Master Thesis, Czech Tech, 2007.

[10] A. Psyllos, C.-N. Anagnostopoulos, and E. Kayafas, "Vehicle model recognition from frontal view image measurements," *Computer Standards & Interfaces*, vol. 33, no. 2, pp. 142–151, 2011.

[11] W. Sun, X. Zhang, S. Shi, J. He, and Y. Jin, "Vehicle type recognition combining global and local features via two-stage classification," *Mathematical Problems in Engineering*, vol. 2017, 2017.

[12] Y. Peng, J. S. Jin, S. Luo, M. Xu, and Y. Cui, "Vehicle type classification using pca with self-clustering," in *2012 IEEE International Conference on Multimedia and Expo Workshops*, 2012, pp. 384–389.

[13] Y. Peng, J. S. Jin, S. Luo, M. Xu, S. Au, Z. Zhang, and Y. Cui, "Vehicle type classification using data mining techniques," in *The era of interactive media*. Springer, 2013, pp. 325–335.

[14] V. S. Petrovic and T. F. Cootes, "Analysis of features for rigid structure vehicle type recognition." in *BMVC*, vol. 2. Kingston University, London, 2004, pp. 587–596.

[15] M.-P. D. Jolly, S. Lakshmanan, and A. K. Jain, "Vehicle segmentation and classification using deformable templates," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 18, no. 3, pp. 293–308, 1996.

[16] K. Kuwabara, Y. Yano, and S. Okuma, "Vehicle appearance model for recognition system considering the change of imaging condition," in *SCIS & ISIS SCIS & ISIS 2008*. Japan Society for Fuzzy Theory and Intelligent Informatics, 2008, pp. 1657–1662.

[17] Y. O. Adu-Gyamfi, S. K. Asare, A. Sharma, and T. Titus, "Automated vehicle recognition with deep convolutional neural networks," *Transportation Research Record*, vol. 2645, no. 1, pp. 113–122, 2017.

[18] J. Wang, H. Zheng, Y. Huang, and X. Ding, "Vehicle type recognition in surveillance images from labeled web-nature data using deep transfer learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 19, no. 9, pp. 2913–2922, 2017.

[19] H. Huttunen, F. S. Yancheshmeh, and K. Chen, "Car type recognition with deep neural networks," in *2016 IEEE intelligent vehicles symposium*

*(IV)*. IEEE, 2016, pp. 1115–1120.

[20] P.-K. Kim and K.-T. Lim, "Vehicle type classification using bagging and convolutional neural network on multi view surveillance image," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 41–46.

[21] X. Ding, J. Wang, C. Dong, and Y. Huang, "Vehicle type recognition from surveillance data based on deep active learning," *IEEE Transactions on Vehicular Technology*, vol. 69, no. 3, pp. 2477–2486, 2020.

[22] J. Taek Lee and Y. Chung, "Deep learning-based vehicle classification using an ensemble of local expert and global networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2017, pp. 47–52.

[23] F. Ullah, H. Anwar, I. Shahzadi, A. Ur Rehman, S. Mehmood, S. Niaz, K. Mahmood Awan, A. Khan, and D. Kwak, "Barrier access control using sensors platform and vehicle license plate characters recognition," *Sensors*, vol. 19, no. 13, p. 3015, 2019.

[24] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *2005 IEEE computer society conference on computer vision and pattern recognition (CVPR'05)*, vol. 1. Ieee, 2005, pp. 886–893.

[25] A. Rio-Alvarez, J. de Andres-Suarez, M. Gonzalez-Rodriguez, D. Fernandez-Lanvin, and B. López Pérez, "Effects of challenging weather and illumination on learning-based license plate detection in noncontrolled environments," *Scientific Programming*, vol. 2019, 2019.

[26] M. Zahedi and S. M. Salehi, "License plate recognition system based on sift features," *Procedia Computer Science*, vol. 3, pp. 998–1002, 2011.

[27] Y. Wang, X. Ban, J. Chen, B. Hu, and X. Yang, "License plate recognition based on sift feature," *Optik*, vol. 126, no. 21, pp. 2895–2901, 2015.

[28] C. Gou, K. Wang, Z. Yu, and H. Xie, "License plate recognition using mser and hog based on elm," in *Proceedings of 2014 IEEE International Conference on Service Operations and Logistics, and Informatics*. IEEE, 2014, pp. 217–221.

[29] B. Li, B. Tian, Q. Yao, and K. Wang, "A vehicle license plate recognition system based on analysis of maximally stable extremal regions," in *Proceedings of 2012 9th IEEE International Conference on Networking, Sensing and Control*. IEEE, 2012, pp. 399–404.

[30] H. W. Lim and Y. H. Tay, "Detection of license plate characters in natural scene with mser and sift unigram classifier," in *2010 IEEE Conference on Sustainable Utilization and Development in Engineering and Technology*. IEEE, 2010, pp. 95–98.

[31] Y. Jamtsho, P. Riyamongkol, and R. Waranusast, "Real-time bhutanese license plate localization using yolo," *ICT Express*, vol. 6, no. 2, pp. 121–124, 2020.

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

[33] M. Onim, S. Hassan, M. I. Akash, M. Haque, and R. I. Hafiz, "Traffic surveillance using vehicle license plate detection and recognition in bangladesh," *arXiv preprint arXiv:2012.02218*, 2020.

[34] R. Smith, "An overview of the tesseract ocr engine," in *ICDAR '07: Proceedings of the Ninth International Conference on Document Analysis and Recognition*. Washington, DC, USA: IEEE Computer Society, 2007, pp. 629–633. [Online]. Available: http://www.google.de/research/pubs/archive/33418.pdf

[35] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2980–2988.

[36] Z. Selmi, M. B. Halima, U. Pal, and M. A. Alimi, "Delp-dar system for license plate detection and recognition," *Pattern Recognition Letters*, vol. 129, pp. 213–223, 2020.

[37] T. Björklund, A. Fiandrotti, M. Annarumma, G. Francini, and E. Magli, "Robust license plate recognition using neural networks trained on synthetic images," *Pattern Recognition*, vol. 93, pp. 134–146, 2019.

[38] X. Yang and X. Wang, "Recognizing license plates in real-time," *arXiv preprint arXiv:1906.04376*, 2019.

[39] R. Laroca, L. A. Zanlorensi, G. R. Gonçalves, E. Todt, W. R. Schwartz, and D. Menotti, "An efficient and layout-independent automatic license plate recognition system based on the yolo detector," *arXiv preprint arXiv:1909.01754*, 2019.

[40] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 4, pp. 834–848, 2017.

[41] J. Zhuang, S. Hou, Z. Wang, and Z.-J. Zha, "Towards human-level license plate recognition," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 306–321.

[42] A. Tourani, A. Shahbahrami, S. Soroori, S. Khazaee, and C. Y. Suen, "A robust deep learning approach for automatic iranian vehicle license plate detection and recognition for surveillance systems," *IEEE Access*, vol. 8, pp. 201 317–201 330, 2020.

[43] H. Li and C. Shen, "Reading car license plates using deep convolutional neural networks and lstms," *arXiv preprint arXiv:1601.05610*, 2016.

[44] L. Zhang, P. Wang, H. Li, Z. Li, C. Shen, and Y. Zhang, "A robust attentional framework for license plate recognition in the wild," *IEEE Transactions on Intelligent Transportation Systems*, 2020.

[45] Tzutalin, "Labelimg," Free Software: MIT License, 2015. [Online]. Available: https://github.com/tzutalin/labelImg