

Research on Multi-Target Vehicle Detection and Tracking Based on Yolo

1st Tchanchou Ngatou Costel

School of Electronic and Information Engineering
Shanghai University of Electric Power
Shanghai, China
E-mail: zhangchucostel@gmail.com

2nd Sanxi Jinag

School of Electronic and Information Engineering
Shanghai University of Electric Power
Shanghai, China
E-mail: samjoe_2018@shiep.edu.cn

Abstract—Vehicle detection and monitoring are gaining importance in traffic management. However, detection is still an issue as vehicles vary in size, which directly affects vehicle counting accuracy. The proposed vehicle detection and counting method first extracts the road surface of the expressway in the image and divides it into distant regions. The newly developed segmentation strategy in the proposed vehicle identification and counting system first extracts the trail's state in the picture furthermore separates it as far as the near areas. This method is important for improving vehicle detection. The above location is then sent to his YOLOv5m network to determine the vehicle type and location. Finally, we validate the proposed methodology using multiple traffic monitoring recordings from different environments. Also, the vehicle detection performance has increased to 99.39% of map compared to YOLOv5 Basic. The research has practical consequences for the management and control of vehicle objects in traffic situations.

Index Terms—traffic management, vehicle detection, vehicle counting, traffic monitoring.

I. INTRODUCTION

Classical machine vision styles and advanced literacy approaches currently fall into two categories. Classical machine vision styles and standard machine vision approaches put excitement into the auto to distinguish it's taken from a static backdrop picture. This approach falls into three orders of magnitude: the background junking method; the non-stop videotape frame differencing method; and the light influx method. Friction is estimated using a videotape frame differencing system based on successive image pixels of two to three videotape forms. Similarly, a threshold isolates the shifting focus regions. Stops can also be detected using this approach, thus reducing noise. Videotape background images are fixed to save background information. Once the backdrop picture on the videotape is established, repeated measures are utilized to build a reference image. The mirrored picture is also linked to a backdrop version to partition shifting items. The light entry shape can represent the stir region in the movie. A deep convolutional network (CNN) was shown. It is very effective in identifying vehicle objects. CNNs are excellent at learning image features and can perform various related tasks such as bracketing and bounding box regression [1]. In general, they have two types of perception styles. The two-layer mode uses multiple styles to build the item search box and uses a convolutional neural network to classify the

items. Although the single-stage mode does not create the search item, it directly translates the SAR image limit issue placement through inverse filtering. The R-CNN (Area) [2] makes overuse of the area look [3] on images in a two-step process. The picture feed to a convent must be of a single node, as the site's inner topology necessitates a long learning curve and massive storage capacity. SPP NET [4] is based on the concept of spatial aggregation matching, which allows the network to accommodate film lands of various sizes while generating stationary waves. RFCN, FPN, and Mask RCNN all feature various advanced convolutional network point generation styles,[5] point selection, and framing capabilities. Chief among the one-step approaches is the Single Shot Multi-Box Sensor (SSD)[6]. SSDs are using Multi Box [7], Region Offer Cross depiction techniques, using networks (RPN), and degraded anchor box sets of various sizes. the rate for better placement of elements. Unlike SSD, the proposed model [8] hub splits the picture into a predetermined number of rasters. Every material monitors predictions such as feature points whose few nodes are inside its boundaries. YOLOv5m [9] included a subcaste BN (batch normalization) that homogenizes the inputs of each subcaste into the network and promotes network confluence. YOLO v5m uses a multiscale training approach to arbitrarily induce, per 10 bits, a new picture width. The network YOLOv5m [10] is used for vehicle object detection.

II. RELATED WORK

A. Engine detection

Intelligent traffic management and road monitoring require vehicle recognition and statistics in road surveillance videos. With my phone, we can collect a large collection of traffic footage for analysis. A broader viewing angle allows you to gauge the road surface more accurately. The size of items within the car varies dramatically at the tilt, or the clarity of identifying small objects far from the road is drastically diminished. When dealing with a complicated video sequence, it is critical to properly diagnose and execute the difficulties described above. This article focuses on the aforementioned issues to uncover potential remedies, and it makes use of vehicle detection data to do so. Small items detected away

from the route have poor sensitivity. When dealing with complicated video sequences, it is critical to examine and apply the aforementioned concerns. In this paper, we concentrate on the aforementioned challenges to get attainable outcomes and apply vehicle detection findings to multi-target detection and vehicle counting challenges. Advanced CNN performs well in object detection. On the other hand, CNNs are extremely delicate to observable alterations in image retrieval [11, 12]. A first-step process predicts an item in a single step. However, the spatial limitations of the grid require a two-step technique that is challenging to implement, particularly for tiny objects. The two-step system groups member search regions into blocks based on set criteria, as well as, if the search area is shorter than the stated criteria, the search regions are interrupted based on the parameters. The supplied parameters are size. As a result, the distinctive structure of little tasks is gone, and the capacity to navigate is diminished. The nature of existence makes no distinction between huge and little objects being in the same order.

B. Engine tracking

Extended vehicle detection tasks, similar to multi-object surveillance, were also important for it. [13] Most advanced multi-images use a detection-based style. The systems use coloring as a model to identify blobs in videotape images prior to recording them. A system that enables scene objects but cannot handle adding new objects or deleting old ones. Multi-object tracking algorithms must consider the similarity of objects within a frame and the related issues of objects between frames. We may utilise tunable cross-correlation to gauge how similar the objects in a frame are to one another. Currently, this problem may be solved by deactivating location detection or location inside the stream. But [14] is using SIFT points to shade the assets to solve the challenges posed by shrinking and lighting moving objects. In this study, we suggest using an algorithm to recognize sphere edges [15]. Spheroids can reach better original sites much more quickly than SIFT. A tiny public data collection for a specific scenario is also available. The capacity of sophisticated neural networks to record changes makes the identification of tiny objects problematic. An established large-scale high-definition vehicle data set can provide a variety of fully annotated vehicle objects from different scenes captured in the traffic domain. The data set is used to assess the effectiveness of a variety of automatic parking methods in terms of vehicle size [16]. The method of recognising dust motes in evidence [17] was used to improve auto-recognition performance.

III. PROPOSED METHOD

The basic structure of the auto detection counting system was clarified in this section. The scene's video data was used as the first input. The road ahead was swept away as well. The YOLOv5m high literacy file detection technique was used to extract mechanical things in the tracking scenario. Finally, spherical point generation was applied to each vehicle block to complete multi-object shading and gather car information. A

pavement segmentation technique, as shown in Fig. 1, was utilised to capture a highway road portion. The route was separated into two halves based on the camera's position. The YOLO v5 m object detection mode can also identify cars on two different sections of the road. This approach may enhance the detection of small items and overcome the problem of recognising delicate objects owing to abrupt changes in object scale. The sphere method was also used to track many objects at the same time. The Sphere method extracts and compares the features of the identified blocks to match correlations between numerous video frames with identical parts. After all, figures are critical computations. The object tracking line was developed to establish the direction of the vehicle and record data such as the number of buses in each order. In terms of video recording, our system is improving object recognition granularity and offers a shadow detection and predictive analytics collecting approach that encompasses the whole field of view [18].

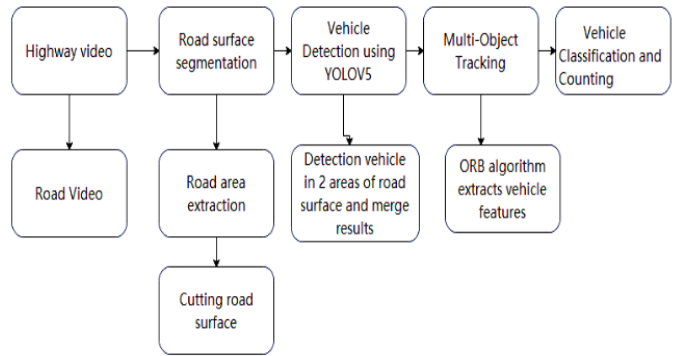


Fig. 1. Overall flow of Method

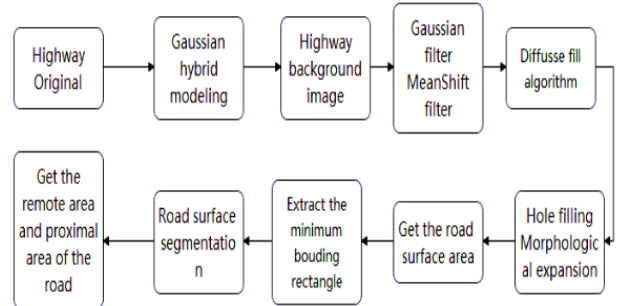


Fig. 2. Average filtration flow of highway

A. Separation of the roadways

This section describes how to extract and segment highway surfaces. Image processing approaches such as Gaussian mixture modeling were used to achieve surface extraction and segmentation [19], enabling excellent vehicle detection results if a fully convolutional descriptor was used methods. The traffic

shots were a broad field of sight. Since this study focuses on cars, the portion of the photo that shows a highway surface was taken. At the same time, the extent of asphalt is focused on one region. The view is based on the camera's angle of view. With this function, we are able to extract the highway lane area in the video. Road surface removal is a process (Fig. 1). To limit the influence of automobiles on road segmentation, a Gaussian mixture modeling technique was employed to extract the backdrop from frame 5.00 at the beginning of the video. The pixel values of an image are Gaussian distributed over a strict mean trust within a concise moment, counting per cell in every shape. If a pixel was off-center, it was considered to be in the foreground. A byte edge was thought to be in scope if its value deviated from the average within a certain range. Mixed Gaussian models were most useful for photos with multiple peak background pixels, such as the road camera on your phone. Photography was used in fieldwork. Following mining, the baseline picture was blurred using a 3*3 core Kalman filter. The trail deck part of a path is chosen by the steppe fill bot as the base, and neighboring areas of the route are filled using the origin point's pixel values. Nearly continuous road surfaces have pixel values that are similar to their original points. Lastly, place To remove the pavement, complete pouring and anatomical expansion were done. Pulling down a few of the truck's fancy freight coverings reveals the aims (Figure 3). (Figure 4).

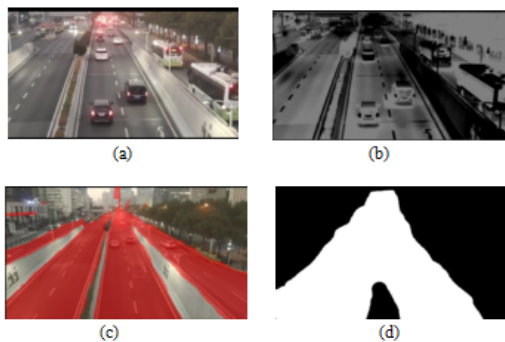


Fig. 3. Procedure for removing the rutting area. a Raw picture; b pic tiffin; c dilute fill; d road surface area mask

B. Vehicle detection using YOLOv5

The YOLOv5 network has the advantage of becoming more important sooner, in contrast to other access pictures in the single-tier. Furthermore, it yielded results similar to previous techniques while maintaining subtlety, and its prediction relied on the regional environment of the source file. Therefore, our proposed model was based on his Yolo fitting from discovery. Anchors in YOLOv5 networks contain many layers to connect. The projects that were carried out at each step can be summarized into three distinct areas of the YOLO-v5 network. For Yolo-v5, the first portion (i.e., called cspdarknet), covers the most common operations in CNNs (complication, movement, max pooling, etc.), as well as future sections. The backbone

concept was common content in several deep literacy networks for object discovery and was used as a simple baseline net. Additionally, the cspdarknet network solves the problem of repeated degrees on large machines. We also significantly reduce pretrained variables and his hanging transactions per second by integrating the change of degree into his area, which improves the speed and sensitivity of inference. The majority of detectors generally do not detect the object as normal. Therefore, we adopt the YOLO v5 model to handle small objects (e.g., Yolo, the focus plane, was the key of Yolo-v5). The focus plane first replicates the shape of the

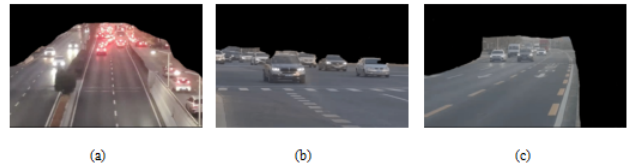


Fig. 4. Results of rutting retrieval for various turnpike instances 1st Drama; 2nd Drama 3rd drama

supplied picture (e.g., 3 x 256 x 256) into 4 copies. The 4 copies were sliced into 4 by scanning with a step size of 2 (that is, 3x128x128). The four sections were also joined depthwise, yielding a 12 x 128 x 128 result, and forwarded into the next layer using controls in 32 cores, yielding a 32 x 128 x 128 n-Output, which was fed to RELU as a stack normalisation and activation function in the sequel layer. Focus subcaste (Fig. 8) Separate the image, transform the spatial information into non-qualitative features, and aggregate the RGB information. It's a fast method that uses GPU math to significantly reduce completion time, but it's designed to transform spatial information into depth information in just a few iterations of image mining. Deep networks were better at anchoring small objects by upsampling the input image instead of down-sampling with YOLO Discovery and inputting the image for certain aspects. (e.g., 416 x 416) Once transmitted to the LAN, because the face was split, the factors of the asphalt in the distance were distorted and increased. As a result, to prevent missing a few object features because the traffic asset was inadequately large, we can identify extra points for small vehicle objects.

For training purposes, a truck image analysis system was developed. The tyre file recognition model can recognise 3 types of trucks: trams, trucks, and swaps (Figure 9). The latter were on the track, We could not find them because there were many bicycles on them. The network receives remote and near-surface parts of the road for discovery. The positions of the vehicle boxes in the two linked regions were plotted inversely to the initial image to capture the appropriate call positions in the initial image. Adopting a car spotting approach to determine the order and position of the vehicles should provide useful records for feature shading. Previous details can be used for vehicle counting, so that vehicle detection systems ignore specific attributes or status of a vehicle. Cars are detected

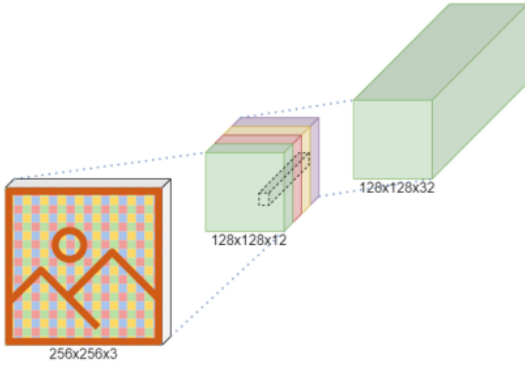


Fig. 5. Basic focus layer of YOLO-v5 model

in the input photos, but the algorithm reclassifies them as buses. YOLO-v5 can accurately categorise features in the input image using the simple YOLO-v5 model, although we still see some false deficiencies (e.g. outside). YOLO-v5 is a deep discovery brand focused on inference with low criteria and an increased frame rate to accommodate reduced outstations. Model changes should not affect this aspect significantly. Such operations usually calculate the average action of the reconnaissance sections on each route.

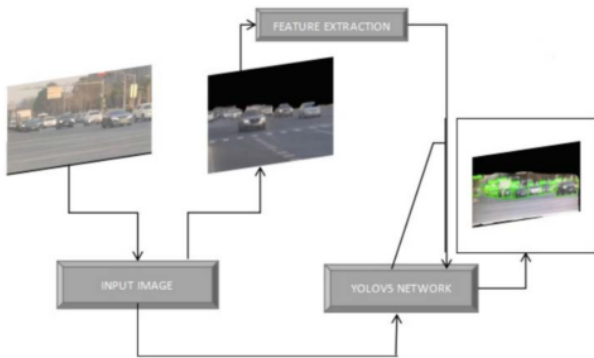


Fig. 6. Structured picture submitted to the detection network, and the detected results are merged. (The color green denote the "car," "bus," and "truck" zones, respective

IV. EXPERIMENTS

A. Dataset

In terms of image access, the picture subset may be split into 03 orders: pics captured by cameras, photos captured by cameras, and captured via means of lenses [20]. A standard data set [21] includes evidence of both lane area as well as general street context, which can be used in independent driving and task problems as well as 3D object detection and shadowing. The Shanghai Canvas business signature data set [22] contains images from cameras covering different light and

precipitation conditions, but no tagged vehicles. This record contains vehicle orders, including the vehicle make, model, and year of manufacture. However, the dataset contains many images. The 28300 printout shows the engine's peak velocity, the count of locks, demotion, and engine sort. 150,200 prints form the complete picture of the vehicle. The Driver File was a good example. [23], which involved 10,000 printouts. This collection classifies vehicles into six instance categories: SUVs, hydrofoils, minivans, trucks, machines, and microbuses. The launch angle, on the other hand, was hopeful, and the motor entity was small. to generalize to CNN training. Introducing the auto recording derived from the viewpoint of the lane safety videotape created by our company. Nevertheless, the dataset contains many images. The printout shows the engine's peak velocity, the count of locks, demotion, and engine sort.

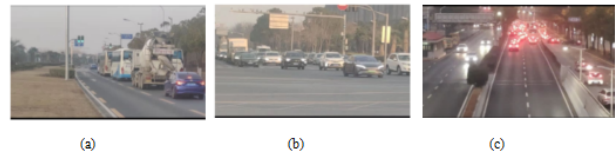


Fig. 7. Scenes captured by mobile phone cameras from multiple highways. Scene 1; b Scene 2. c scene 3

The print indicates the overall appearance of the vehicle., Surveillance cameras recorded recordings, and the BIT vehicle dataset [24], including printouts, was shown. This dataset classifies vehicles into six categories: SUVs, hydrofoils, minivans, trucks, machines, and microbuses. The launch angle, on the other hand, was hopeful, and the motor entity was too small to generalize to CNN training. We will introduce the vehicle image seen from the lane monitoring videotape manufactured by our company. The dataset images were from his Lingang and YangShupu street videos in China (Fig. 7). Mobile phone cameras were installed on the side of the road, and we are installing them in 2 locations for lane monitoring. It has no predefined adaptive wide-angle positions. Pictures from this vantage point wrap around long stretches of road and contain automobiles of all shapes and sizes. The images were recorded in the source by a mobile phone camera, in various images, on several occasions, in various indoor situations. Every retrieval classifies lorries into 3 parts: commuters, passenger cars, and transit connections.



Fig. 8. Vehicle labeling category of the dataset

V. TRACKING OF SEVERAL OBJECTS

This section shows how to track many things using object boxes found in the part on vision-based based employing YOLOv5. The Arrow technique was successfully investigated to extract the properties of the identified autos. The Shard classifier beats other schemes in terms of processor speed and pairing charge. This model was a viable alternative to the Slog and Wave image analysis methodologies. The image approach employs the driver to find sides and attributes as from the chart provided by the brief example to find a descriptor. Having received the images, the descriptors were computed using the Brief method after the feature points were obtained. The main service was the center of the circle, as well as the midpoint of a point zone, which serves as the cross of the set of points. Since the shot may be flipped, it is possible to rotate the grid as well. As a result, the minutiae descriptors were rotationally consistent. Even if the angle of view changes, you can still propose a fixed point. When the number of matching points collected exceeds the defined threshold, the points are considered successfully matched and the object's matching field is drawn. The RANSAC algorithm was utilized for point recovery, which allows incorrect input areas to be eliminated from appropriate calculations, as well as a high forecast.

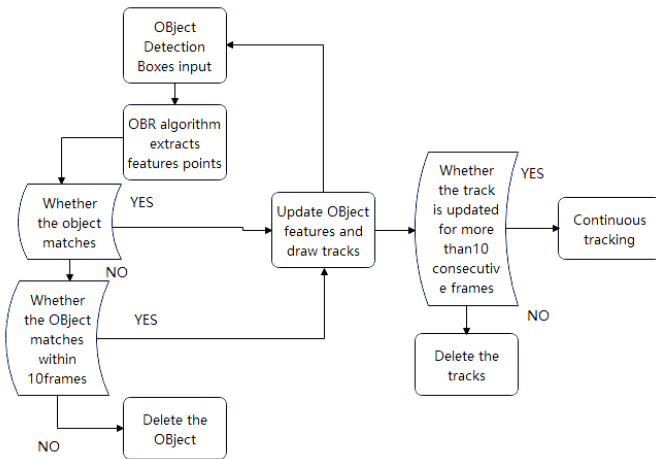


Fig. 9. Process of multi-object tracking

Concordance in both frames (Fig. Unless the vaccination zone and the sensor crate in the following quarter are within the center-to-center distance, criterion, Fig. 11)

$$T = \frac{\text{box height}}{0.25} \quad (1)$$

This relates to the highway under examination as seen from a wide-angle camera. The lane layer reported by the gadget in this photo was well below. The track descends for ten consecutive frames of video. So, if we don't stiffen the images by around 10 frames, they will break. The tracking seen on the whole trajectory monitoring videotape and the final object identification result are modified by processing the information

provided above because the auto route and device wire overlap till the beginning.

A. Interpretation of the findings

This section describes appraisal of the styles stated in the Styles portion. In the wheel record portion, we attempted to a van item record. In while experiments, We were using detailed tracking scenes in 3 different regions, as shown in Table 1.

TABLE 1. Number of objects under different detection methods

| Scenes | Total number of vehicle objects | Total number of images | full-image detection method | Actual number of vehicles | Actual number of vehicles | Our method | | |
|---------|---------------------------------|------------------------|-----------------------------|---------------------------|---------------------------|------------|---------|-------------|
| | | | | | | video | vehicle | Remote area |
| Scene 1 | Car | 3,000 | 6,128 | 8,430 | 493 | 6,616 | 6,849 | 8,550 |
| | Bus | 3,000 | 535 | 459 | 92 | 379 | 582 | 483 |
| | Truck | 3,000 | 5,311 | 5,320 | 840 | 4,703 | 5,792 | 5,471 |
| Scene 2 | Car | 3,000 | 1,843 | 3,615 | 192 | 3,356 | 1,914 | 3,654 |
| | Bus | 3,000 | 194 | 364 | 82 | 295 | 207 | 382 |
| | Truck | 3,000 | 3,947 | 4,709 | 922 | 3,738 | 4,169 | 4,731 |
| Scene 3 | Car | 3,000 | 1,774 | 2,336 | 224 | 2,188 | 1,834 | 2,352 |
| | Bus | 3,000 | 415 | 516 | 56 | 495 | 483 | 529 |
| | Truck | 3,000 | 3,678 | 3,490 | 731 | 2,662 | 3,726 | 3,507 |

B. Detecting vehicles and traffic practice

For the identification of the vehicle objects and the recording of the data for the entire training, we used the YOLOv5 network. For the division of the data set in the network training, there is no ideal result. Our strategy of splitting the records and dividing the data follows the practice. To divide the data set into a training set with a value of 60% and a test set with a value of 30%. The photos for the training and test sets were randomly selected from the dataset, which includes prints. The test and training sets were sufficient to obtain the model because the dataset contains many photos. To produce an accurate model, the pace of the training set needs to be quickly photographed, which guides us through several vehicle instances to make realistic representations of buses, trucks, and truck destinations. These were included in the training set. The result set contains 2300 impressions of vehicle stimuli that were completely separate from the data set, which was sufficient to evaluate the figure. We set the mass damping to 0.0008 or the closure value to 0.9 for a limit of 60,500 duplicated exercises. The rate for the first 30,000 copies was 0.01.

After 30,000 duplicates, the rate was reduced to 0.001. This system has significantly reduced losses. A k-means approach was used to modify the anchor field to make it more suitable for annotation of the recorded field. On the training dataset, we estimated the size of the set of the nodes at a routing specificity of 5 records with an average IOU of 95.5 percent: (15,357,33,542); (188, 17, 37); (143,60,63); (88,73,105,53); (133,95,18). We did not remove samples smaller than 1 pixel to improve object detection. When we configured the network's image input, the network resolution became 832*832 instead of 416*416. If it is a Yolo sub-casting problem, after increasing the input resolution, the network resolution will also increase,

which will improve the object detection accuracy. See Figure 11.



Fig. 10. Result of detecting video objects by frame. Each green cell is labeled with a region for cars, buses, and trucks. our way. b Frame detection method

Increasing the input resolution and then outputting the network with a Yolo layer will increase the resolution proportionally and boost object performance. The trained model is used to detect vehicles under various road conditions using a 3000-frame continuous image sequence. Road areas are captured and isolated before inclusion in the vehicle detection mesh.



Fig. 11. Vehicle and detection datasets using fine-tuned YOLO-v5

TABLE 2. Comparison of actual number of vehicles by different methods

| Vehicle | Remote area | Proximal area | | Average correct rate | | | |
|---------------------------|------------------|------------------|------------------|----------------------|------------------|--------|--------|
| | | Full-image | Our method | Full-image | Our method | | |
| Our method | | | | | | | |
| | | | | | | | |
| Category | detection method | detection method | detection method | detection method | detection method | | |
| | Car | 98.95% | 12.58% | 85.54% | 98.54% | 97.45% | 48.06% |
| Actual number of vehicles | Bus | 89.95% | 20.38% | 97.15% | 73.86% | 93.77% | 54.33% |
| | Truck | 95.51% | 22.31% | 98.01% | 84.11% | 98.16% | 51.21% |
| Overall correct rate | | 92.14% | 16.26% | 98.32% | 88.28% | 97.36% | 50.47% |

A method for identifying photos having a 1920 x 1080 resolution in a network (without road surface segmentation) Table 2 and Figure 11 illustrate the results. Table 2 compares the number of object detections made using various approaches to the number of actual cars. in comparison to reality If the short-range objects on the road are really large, our method approximates the actual number of automobiles. Even when the item is small and far away from the road, the observed

deviation was less than 10%,the Two examples are inference speed, which is frequently connected to frames per second (FPS), and the number of features, which is usually a great measure of style complexity.

C. Evaluation Metrics

We utilised a number of criteria [0, 0.1, 0.2, and 1] for the maximum recall Pmax, which is bigger than any parameter (the experimental limit is 0.25) (recall). The accuracy was determined, and AP is the average of these P maximums (recall). This number is used to signify the model's quality. Several criteria are used to determine the efficacy of machine-learning methods. Precision (P) is the proportion of genuine positives discovered between every one identified.

$$ap = \frac{1}{11} \sum_{recall=0} (pmax_{recall}), recall \in [0, 0.1, 0.2..] \quad (2)$$

$$map = \frac{\sum ap}{class\ number} \quad (3)$$

The correctness and recall assessments are as below.:

$$Correctness = \frac{TA}{TA + FA} \quad (4)$$

$$Recall = \frac{TA}{TA + FN} \quad (5)$$

where TA, FN, and FA imply the number of true defects and false defects. The final result in the chart is 99.39 % , demonstrating that the strategy of finding and classifying colorful car details is effective. The above analysis shows that the average throughput is 92.64 % . It demonstrates proper placement and classification of various vehicle objects, as well as improved detection results for multiple objects.The average complete of all results with the intersection of 50 %and 95% units (IOU) is represented by mAP.5 and mAP.95. where IOU is the result of sending the intersection of the input image (found highly valued and base value) constrained by the unity of objects . An average accuracy (AP) is also calculated based on detection of specific classes with an IOU of less than 50 % or 95 % . At last, We calculated the ordinary accuracy (map) by using surface over all categories.

D. Monitoring and tagging

Following the acquisition of the object boxes, we used the ORB feature point matching technique to perform vehicle tracking and trajectory analysis. Each asset's related projected ORB position in the experiment is produced by a match score inferior to 10. The tracking trajectory was determined using the acquisition lines. Table 3

We did a test with 3 other movies. This is the same situation as in the auto screening and server skill sections. We evaluated the speed of the system proposed in this study using the reasonable speed, which is measured by the days allotted for the method to handle the video and the time it takes to play the original video. In the formula In Figure 4, network lag is

TABLE 3. Number of objects with different detection methods

| Scenes Video Frames Vehicle category | Scene 1 | | | Scene 2 | | | Scene 3 | | | Direction correct rate | |
|---|---------|-------|-------|---------|-------|-------|---------|-------|-------|---------------------------|-------|
| | A | | | B | | | C | | | | |
| | Car | Bus | Truck | Car | Bus | Truck | Car | Bus | Truck | | |
| Our method | 19 | 11 | 3 | 110 | 30 | 21 | 185 | 121 | 22 | | |
| Actual number of vehicles | 21 | 11 | 3 | 117 | 33 | 22 | 197 | 130 | 24 | | |
| Direction A | | | | | | | | | | | |
| Extra Number | 2 | 0 | 0 | 7 | 3 | 1 | 12 | 9 | 2 | 0.95 | |
| Missing number | 0 | 0 | 0 | 1 | 0 | 1 | 2 | 4 | 1 | | |
| Correct rate | 0.901 | 1 | 1 | 0.95 | 0.93 | 0.85 | 0.91 | 0.89 | 0.833 | | |
| Our method | 51 | 47 | 4 | 107 | 89 | 17 | 300 | 184 | 15 | | |
| Actual number of vehicles | 53 | 48 | 4 | 115 | 97 | 17 | 311 | 192 | 17 | | |
| Direction B | | | | | | | | | | | |
| Extra Number | 2 | 1 | 0 | 8 | 8 | 0 | 11 | 8 | 2 | 0.94 | |
| Missing number | 0 | 1 | 0 | 3 | 2 | 0 | 4 | 4 | 0 | | |
| Correct rate | 0.963 | 0.947 | 1 | 0.888 | 1.101 | 1 | 0.939 | 0.93 | 0.892 | | |
| Real time rate | | 1.95 | | | 1.1 | | | 1.2 | | | |
| Average correct rate | | 0.945 | | | 0.922 | | | 0.917 | | | 0.935 |

the optimal time to process the video, as well as the timing of video is the time it takes to play the movie. The lower the implied volume number, the slower the reliable computations. If the real value is null and the speed is less than or equal to one, the visual input can be handled in real time.

$$\text{Real time rate} = \frac{\text{running time of the process}}{\text{duration of the video}} \quad (6)$$

According to the data, the average accuracy for direction finding and vehicle counting is 92.46%, 92.52%, and 98.61%, respectively. Vehicle ratings in road surveillance videos are displayed as small specks and are easily blocked by large vehicles. Since multiple vehicles were moving at the same time, it affected the accuracy of distance measurement. The original video has a frame rate of 30 fps. Velocity calculations show that the vehicle tracking method based on the ORB function is faster. The speed of the system is related to the type of cars by scene. This increases system processing time. The vehicle counting method proposed in this research is generally very close to real-time processing.

TABLE .4 The network model’s specificity

| Parameters | Car | Bus | Truck | Precision | Reccal all | Average IOU | mAP |
|------------|--------|--------|--------|-----------|------------|-------------|--------|
| Results | 92.46% | 92.57% | 98.61% | 0.833 | 0.9886 | 95.00% | 99.39% |

VI. CONCLUSION

In this work, we constructed a relevant auto asset collection by the camera’s position proposed an entity recognition and shadowing algorithm to lane video sequences captured by phones. The lane roadside area baseline yielded yet another acceptable ROI area. On the basis of the collection of annotated lane vehicle objects, the YOLOv5 object identification algorithm created a whole chain of lane engine detection brands. The test results demonstrated that the suggested vehicle recognition and strategy for tracking videotape sequences

captured by mobile phones work well and are viable. Using the multiscale coefficient attention mechanism, the updated interpretation exceeds his YOLO-v5 for this specific operation, with an accuracy of over 0.833 compared to 0.67 for the YOLO-v5 base style increase. The suggested model also marginally enhanced recall and mAP values while keeping the same number of biddable variables. However, the suggested detection and tracking mechanism should be improved. Faster car detectors tailored particularly for traffic situations are envisaged in future development.

REFERENCES

- [1] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, “Object detection with deep learning: A review,” *IEEE transactions on neural networks and learning systems*, vol. 30, no. 11, pp. 3212–3232, 2019.
- [2] R. Girshick, J. Donahue, T. Darrell, and J. Malik, “Rich feature hierarchies for accurate object detection and semantic segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
- [3] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders, “Selective search for object recognition,” *International journal of computer vision*, vol. 104, no. 2, pp. 154–171, 2013.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, “Spatial pyramid pooling in deep convolutional networks for visual recognition,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 9, pp. 1904–1916, 2015.
- [5] S. Ren, K. He, R. Girshick, and J. Sun, “Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 6, pp. 1137–1149, 2017, doi: 10.1109/TPAMI.2016.2577031.
- [6] Y. O. L. Once, “Unified, real-time object detection/redmon j., divvala s., girshick r., farhadi a,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [7] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, “Going deeper with convolutions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [8] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester, “Cascade object detection with deformable part models,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, San Francisco, 2010*, pp. 2241–2248, doi: 10.1109/CVPR.2010.5539906
- [9] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, “A unified multi-scale deep convolutional neural network for fast object detection,” in *European conference on computer vision*. Springer, 2016, pp. 354–370.
- [10] X. Hu, X. Xu, Y. Xiao, H. Chen, S. He, J. Qin, and P.-A. Heng, “Sinet: A scale-insensitive convolutional neural network for fast vehicle detection,” *IEEE transactions on intelligent transportation systems*, vol. 20, no. 3, pp. 1010–1019, 2018.
- [11] A. A. Nielsen, “The regularized iteratively reweighted mad method for change detection in multi-and hyperspectral data,” *IEEE Transactions on Image processing*, vol. 16, no. 2, pp. 463–478, 2007.
- [12] D. Rosenbaum, F. Kurz, U. Thomas, S. Suri, and P. Reinartz, “Towards automatic near real-time traffic monitoring with an airborne wide angle camera system,” *European Transport Research Review*, vol. 1, no. 1, pp. 11–21, 2009.
- [13] J. Canny, “A computational approach to edge detection,” *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986.
- [14] P. Negri, X. Clady, S. M. Hanif, and L. Prevost, “A cascade of boosted generative and discriminative classifiers for vehicle detection,” *EURASIP Journal on Advances in Signal Processing*, vol. 2008, pp. 1–12, 2008.
- [15] Q. Fan, L. Brown, and J. Smith, “A closer look at faster r-cnn for vehicle detection,” in *2016 IEEE intelligent vehicles symposium (IV)*. IEEE, 2016, pp. 124–129.
- [16] H. Asaidi, A. Aarab, and M. Bellouki, “Shadow elimination and vehicles classification approaches in traffic video surveillance context,” *Journal of Visual Languages Computing*, vol. 25, no. 4, pp. 333–345, 2014.
- [17] Q.-L. Li and J.-F. He, “Vehicles detection based on three-frame-difference method and cross-entropy threshold method,” *Computer Engineering*, vol. 37, no. 4, pp. 172–174, 2011.

- [18] Z. Zhu, D. Liang, S. Zhang, X. Huang, B. Li, and S. Hu, "Traffic-sign detection and classification in the wild," in Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 2110–2118.
- [19] T. Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal Loss for Dense Object Detection," Proc. IEEE Int. Conf. Comput. Vis., vol. 2017-October, pp. 2999–3007, 2017, doi: 10.1109/ICCV.2017.324..
- [20] A. Rangesh and M. M. Trivedi, "No Blind Spots: Full-Surround Multi-Object Tracking for Autonomous Vehicles Using Cameras and LiDARs," IEEE Trans. Intell. Veh., vol. 4, no. 4, pp. 588–599, 2019, doi: 10.1109/TIV.2019.2938110.
- [21] J. Xing, H. Ai, and S. Lao, "Multi-object tracking through occlusions by local tracklets filtering and global tracklets association with detection responses," in 2009 IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2009, pp. 1200–1207.
- [22] . Zhou, Y. Yuan, and C. Shi, "Object tracking using sift features and mean shift," Computer vision and image understanding, vol. 113, no. 3, pp. 345–352, 2009.
- [23] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in 2011 International conference on computer vision. Ieee, 2011, pp. 2564–2571.
- [24] P. Soviany and R. T. Ionescu, "Optimizing the trade-off between single-stage and two-stage deep object detectors using image difficulty prediction," Proc. - 2018 20th Int. Symp. Symb. Numer. Algorithms Sci. Comput. SYNASC 2018, pp. 209–214, 2018, doi: 10.1109/SYNASC.2018.00041.
- [25] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for, autonomous driving? the kitti vision benchmark suite. in, 2012," in IEEE Conference on Computer Vision and Pattern, Recognition, pp. 3354–3361.