





VSNet: Vehicle State Classification for Drone Image with Mosaic Augmentation and Soft-Label Assignment

Youlkyeong Lee , Jehwan Choi , and Kanghyun Jo 

Department of Electrical, Electronic and Computer Engineering, University of Ulsan,
Ulsan 44610, South Korea
choijh1897@gmail.com, acejo@ulsan.ac.kr

Abstract. Numerous architectures are under development to comprehend object information and background in images by analyzing features extracted through Convolutional Neural Networks (CNNs). Autonomous driving requires understanding diverse information and collecting data from heterogeneous environments to generalize classification models. However, the patterns of feature maps extracted through convolution layers in drone image data, which encompass assorted types of vehicles and road shapes, tend to be simple, leading to overfitting during model training. To prevent overfitting, this study applies Mosaic Augmentation to increase data diversity and brings generalization to the data. This data augmentation method randomly combines four selected images to create a new mosaic image. Soft-label Assignment is used to determine the labels of the mosaic images. The dataset is collected using a drone flying along roads, and approximately 4,000 images are used for training. In the experiment, the classification performance of vehicle status is listed based on the weight of the loss function of the soft label and hard label. Having achieved an accuracy of 83.41%, the effectiveness of the proposed method is compared with dilated residual networks in terms of improving model performance.

Keywords: Drone image · Transportation system · Vehicle state · Classification

1 Introduction

Generating desired information through algorithms using various CNN models is crucial for collecting and analyzing image-based traffic information. Utilizing drones for data analysis provides a different perspective on the images compared to the black box or installed cameras in conventional vehicles. Therefore, a drone capturing a wider area at once is essential for comprehending comprehensive

This work was supported by the National Research Foundation of Korea (NRF) grant funded by the government (MSIT) (No. 2020R1A2C200897212).

© The Author(s), under exclusive license to Springer Nature Singapore Pte Ltd. 2023
N. T. Nguyen et al. (Eds.): ACIIDS 2023, LNAI 13995, pp. 109–120, 2023.
https://doi.org/10.1007/978-981-99-5834-4_9

vehicle movement and flow. Understanding vehicle movement is necessary and can be categorized into three basic states: normal state, where the vehicle moves along the lane; lane-changing state, where the vehicle changes its lane; and stop state, where the vehicle comes to a halt.

This research aims to design a CNN model that classifies three vehicle movement states through drone images and applies mosaic data augmentation [1] and soft label assignment [2]. The dataset is collected data using a drone, which captured images with a bird's-eye view of the road. This work utilizes mosaic augmentation to increase data diversity and prevent overfitting during model training. This data augmentation technique randomly combines four selected images to generate a new mosaic image. Soft-label Assignment is used to determine the labels of the mosaic images. These techniques demonstrate the potential of drones for traffic information analysis and the effectiveness of the proposed methodology in improving classification model performance for autonomous driving systems.

2 Related Work

2.1 Autonomous Vehicle Dataset for Object Detection

Images and annotation data collected in the past environment on various roads have been continuously accumulated. With the advancement of object classification and detection technology, autonomous driving technology is being developed rapidly. The Cityscapes [3] and KITTI [4] datasets were created as datasets for autonomous vehicle research [5–10]. The dataset generated image collection and annotation data for traffic conditions on the road through cameras installed in the vehicle. The KITTI dataset also includes 3D bounding box location and camera calibration information through a 3D laser scanner.

2.2 Drone-Based Dataset for Object Detection

Stanford drone dataset [11] is the first public aerial image dataset using drones. This dataset contains ten kinds of tracking information (Track ID, (xmin, ymin), (xmax, ymax), frame, lost, occluded, generated, label) about objects on the road in the video image. Images taken at eight locations on the Stanford campus were collected. The targets are six classes (Bicyclist, Pedestrian, Skateboarder, Cart, Car, and Bus). However, the annotation quality of the bounding box of the object is roughly expressed, which has a problem with the performance of the object detection algorithm. The VisDrone [12] dataset is a large-scale drone image produced by AISKEYEYE team at Lab of Machine Learning and Data Mining, Tianjin University, China. The dataset aims to develop applications that can be used for computer vision through drones. Through cameras installed in drones, 288 video images were collected from 14 urban areas in China. It produced 2.6 million bounding boxes, including ten classes (pedestrian, person, bus, car, van, truck, bicycle, awning tricycles, motorcycles, and tricycles). Data validation is

tested through VisDrone challenge [13–15] and various kinds of research [16–18] are utilized. The Institut für Kraftfahrwesen Aachen research team had built a drone-based road user trajectory dataset for various situations. The test for vehicles related to autonomous driving is conducted based on the scenario. Therefore, we present reliable and high-quality data criteria. The highD [19] is a large-scale vehicle trajectory dataset for German high roads. It includes six locations, 16.5 h, and 110,000 trajectory information. In inD [20], automated vehicles require data-based analysis methods to understand complex environments. By collecting road images using drones, it was proposed to collect road trajectories and natural road conditions through vehicle movement. Finally, the dataset provides a dataset including road conditions and vehicles, bicycles, and pedestrians over four kinds of German intersections. The roundD [21] includes the movement trajectories of cars, vans, trucks, buses, pedestrians, bicycles, and motorcycles in three traffic circles in Germany. In addition, positions, headings, speeds, accelerations, and classes of objects were extracted from the video and provided as data.

3 Proposed Algorithm

Figure 1 illustrates the overall process for classifying vehicle status. The process consists of four components: 1) Vehicle detection with YOLOv5, 2) Mosaic data augmentation, 3) Soft-label Assignment, and 4) a network for vehicle state classification. These four components work together to determine the movement status of vehicles. This chapter provides an explanation of the proposed methods in detail.

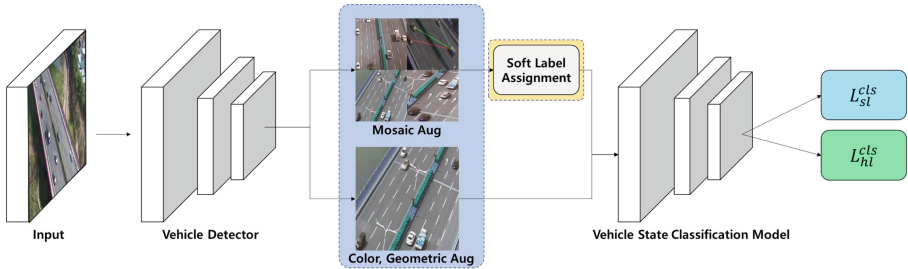


Fig. 1. An overview process of vehicle state classification that contains object detection, data augmentation (mosaic augmentation), soft label assignment, and VSNet (vehicle state network).

3.1 Vehicle Detection

This study first presents an approach for detecting vehicles using YOLOv5 [22], an advanced object detection algorithm that has achieved state-of-the-art performance on a variety of visual recognition tasks. YOLOv5 is an abbreviation

for “You Only Look Once version 5”, and is an extension of the original YOLO algorithm with improvements in speed and accuracy. It is based on a deep neural network architecture that efficiently extracts features from images and predicts object bounding boxes and class probabilities in a single forward pass.

The YOLOv5 algorithm comprises two main components: a feature extraction backbone and a detection head. The backbone network is built on efficient architecture, which has been shown to be highly efficient and effective in a wide range of vision tasks. The detection head employs anchor boxes and grid cells to predict object locations and classes at multiple scales. To adapt YOLOv5 for vehicle detection task, the train fine-tunes the model on a custom dataset of drone flight images using transfer learning. Vehicle types are limited to car, truck, and bus. Specifically, this work initializes the network with pre-trained weights on the COCO dataset [23] with advanced data augmentation and optimization techniques.

3.2 Mosaic Data Augmentation

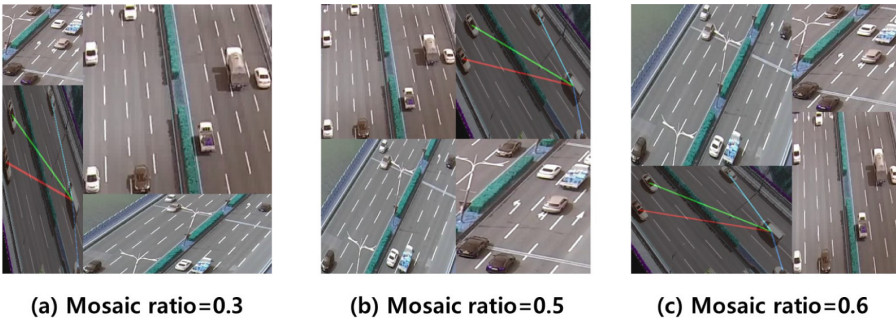


Fig. 2. With the several mosaic ratios, mosaic augmentation generates mixed 4 images to a new image.

Mosaic data augmentation is proposed in YOLOv4 [1] as a technique for augmenting data. This method involves selecting four images from the dataset and arranging them in a manner that is determined by the mosaic ratio, denoted by \mathcal{M}_r so that they are represented as a single image in Fig. 2. \mathcal{M}_r is chosen randomly from the range [0.3, 0.7]. $I(i)$ is image among dataset at index i . Index, i is randomly selected. Every $I(i)$ is resized to 512×512 . w_i and h_i denote the width and height of $I(i)$. $I(n)$ contains 4 images that n is an order, $n = 0, \dots, 3$. Based on \mathcal{M}_r , the width and height sizes of each of the four images are determined as follows:

$$I(n) = \begin{cases} \text{new } w_0 = \mathcal{M}_r \times 512, \text{ new } h_0 = \mathcal{M}_r \times 512 & \text{if } n \text{ is } 0 \\ \text{new } w_1 = 512 - \text{new } w_0, \text{ new } h_1 = 512 - \text{new } h_0 & \text{if } n \text{ is } 1 \\ \text{new } w_2 = \mathcal{M}_r \times 512, \text{ new } h_2 = 512 - \text{new } h_0 & \text{if } n \text{ is } 2 \\ \text{new } w_3 = 512 - \text{new } w_2, \text{ new } h_3 = \mathcal{M}_r \times 512 & \text{if } n \text{ is } 3 \end{cases} \quad (1)$$

This approach enables the model to learn from multiple images simultaneously, improving generalization by incorporating diverse contextual information into a single image. Mosaic is applied by selecting random numbers in quantity equal to the batch size for each iteration.

3.3 Soft Label Assignment

After applying mosaic augmentation in this study, a method of soft label assignment for label allocation is proposed. In the mosaic image, four original images correspond to four labels. A soft label is created by referring to label smoothing [2]. The soft label, $\mathcal{S}(x)$ is shown in Eq. (2). In training sample x , $h_i(k|x)$ represents the hard label distribution of the four images at classes, $k \in 0, 1, 2$ and index of distribution i . The hyperparameter α is assigned a weight value between 0 and 1. The value of K , which denotes the number of images, is 4. Equation (2) are defined as the following:

$$\mathcal{S}(x) = \frac{1}{N} \sum_{i=1}^N \{(1 - \alpha)h_i(k|x) + \alpha/K\} \quad (2)$$

The ground truth label distribution multiplies the weight α and interpolates the hard label through the α/K . The label of the mosaic image generates a soft label by calculating the average value from $(1 - \alpha)h_i(k|x) + \alpha/K$. Equation (2) represents the ground truth soft label, which adjusts the ground truth label distribution by applying label smoothing to mosaic images for classification models. In conclusion, these methods adapt the use of mosaic augmentation and soft label assignment resulting in improved classification model performance.

3.4 Vehicle State Classification

This classification model is adapted [24] as the previous work. The proposed model comprises the Wide Area Feature Extraction (WAFE) module and Deformable Residual (DR) module. These modules play critical roles in extracting and focusing on feature information. The following section provides a detailed layer-by-layer explanation of these modules.

Wide Area Feature Extraction Module (WAFE module). To classify the state of the target vehicle, the input image considers the position and state of the surrounding vehicles. Figure 5 shows that the vehicles in the image are mostly separated. To exclude unnecessary information like background, the first convolutional layer passes a 5×5 kernel size with 64 filters, a stride of 4, and a dilated ratio of 3. Next 1×1 convolutional layer extends the number of channels, 64 to 128. To stabilize the learning process on the feature map, batch normalization (BN) [25] is performed after all convolutional layers, and the proposed network employs Gaussian Error Linear Units (GELU) [26] as the activation function. The feature map is further processed by dividing the 32 channels into

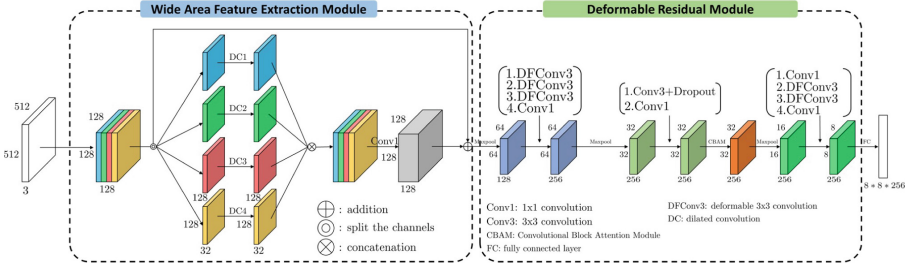


Fig. 3. The left module in the illustration is the WAFE module. It uses a dilated convolutional layer to appropriately extract features while reducing computation when objects in the image are far apart. The right module is the DR module, which is designed to extract meaningful features for vehicle status judgment by applying a variety of receptive fields using a deformable convolutional layer.

four groups, and each group is passed through four kinds of dilated ratio, [1, 3, 5, 7] of 3×3 convolutional layer. The four groups of outputs are concatenated, and a 1×1 kernel is applied. Additionally, a residual block is used to incorporate previous information before the maxpooling operation into the feature map.

Deformable Residual Module (DR Module). Deformable residual is modified from deformable convolutional layer [27] to extract flexible spatial information through output feature from WAFE module. As illustrated in Fig. 4(a), the traditional 3×3 convolutional layer has a fixed receptive field in the image area, represented by the red and blue dots. However, in the image data used for vehicle detection, the vehicles are often separated from each other. Therefore, using a fixed receptive field would extract feature information that includes unnecessary background information. To address this issue and perform more effective convolutional operations, deformable convolution is employed. Figure 4(b) shows how deformable convolution generates an offset as the convolutional layer and performs convolution operations through the offset information.

$$\mathbf{y}(\mathbf{p}_0) = \sum_{\mathbf{p}_n \in \mathcal{R}} \mathbf{w}(\mathbf{p}_n) \cdot \mathbf{x}(\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n) \quad (3)$$

Equation (3) represents the offset that determines the kernel position in deformable convolution. \mathbf{y} represents the output feature map. The kernel grid, \mathcal{R} , is defined as the receptive field, where $\mathcal{R} = (-1, -1), (-1, 0), \dots, (0, 1), (1, 1)$. The convolution occurs at the pixel position of the input image \mathbf{x} , which is \mathbf{p}_0 , and at the individual positions in \mathcal{R} , which are \mathbf{p}_n , along with the offset, $\Delta \mathbf{p}_n$. In particular, the offset $\Delta \mathbf{p}_n$ value is generated based on the convolution layer value and is trained in each iteration. Thus, $\mathbf{p}_0 + \mathbf{p}_n + \Delta \mathbf{p}_n$ ultimately determines the position of the input value, and the convolution operation is performed by multiplying it with the convolution kernel weight, $\mathbf{w}(\mathbf{p}_n)$, at that position.

The Convolutional Block Attention Module (CBAM) [28] allows for complementary attention of both channel-wise and spatial-wise information and is

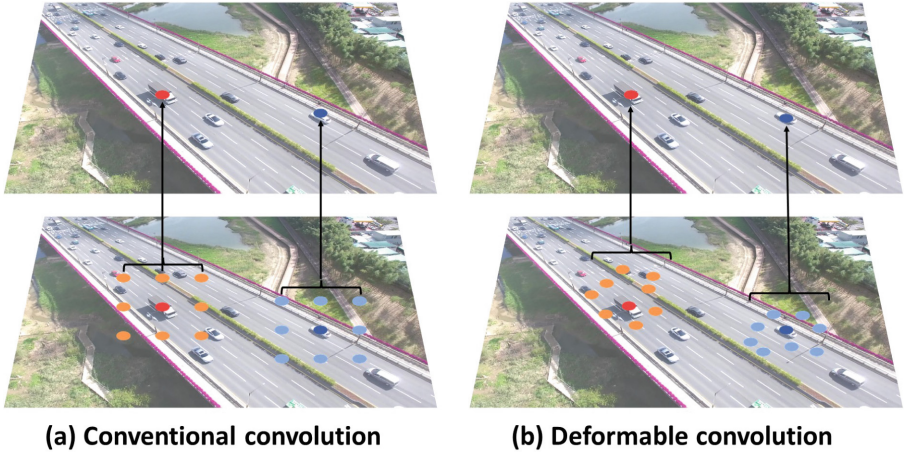


Fig. 4. (a) 3×3 conventional convolution, (b) 3×3 deformable convolution, In deformable convolution, deep red and dark blue dots are focused on the vehicle in the input image. (Color figure online)

applied through the output of three deformable convolutional layer operations. The fully connected layer receives the feature map that has been calculated by two deformable convolutional layers.

Loss Function. During training, hard and soft labels are utilized to adjust the loss function. For hard labels, the original loss function uses in this study is Focal Loss [29], which helps to balance the training process and prevent bias towards one class when dealing with data imbalance issues. For soft labels, the loss function is the mean squared error (MSE). Since the value of the soft labels is a float number, it has been computed average value. The proposed total loss function, \mathcal{L} is as follows:

$$\mathcal{L} = \alpha_l \mathcal{L}_{hl}^{cls} + (1 - \alpha_l) \mathcal{L}_{sl}^{cls} \quad (4)$$

\mathcal{L}_{hl}^{cls} is the loss calculated for the hard label, and \mathcal{L}_{sl}^{cls} represents the loss result for the soft label. The parameter, α_l assigns weights to the hard and soft losses. Since soft labels are selected less frequently than hard labels, the weight of hard labels is higher. α_l is 0.9.

4 Experiment

Drone Image Dataset: The dataset for drone images is captured from a top-down perspective in Fig. 5, and vehicle detection is performed using YOLOv5 large model on the collected images. After detection, the image is cropped based on the five vehicles surrounding the target vehicle. The dataset consists of three classes: lane_change, safe, and stop, and the total number of training and test data is shown in Table 1.

Configuration Details. In this study, the Adam optimizer [30] is employed and the learning rate is set to 0.001. Epoch is 200. Four NVIDIA RTX 3090 GPU, each with 24GB of memory, are used, with a batch size of 16.

Table 1. Information of train and test dataset for vehicle state classification.

Class	train	test	Total
lane_change	860	214	1,074
safe	1,241	310	1,551
stop	1,222	305	1,525
Total	3,323	829	4,152



Fig. 5. Illustration for drone image dataset. The view in the picture is bird's-eye view.

Object Detection. In this study, YOLOv5 [22] is adopted as the object detection algorithm, and car and truck are the two classes considered for train and test. Table 3 presents the object detection performance for the train and test datasets. The training on 9,776 images a performance of 95.75 mAP(AP_{50}) and 83.8 mAP($AP_{50:95}$), while testing on 2,200 images a performance of 91.8 mAP(AP_{50}) and 80.3 mAP($AP_{50:95}$). Utilizing this detector, other traffic videos are analyzed to identify and extract vehicle information, including their position and class (Table 2).

VSNet Performance. The performance of the network for classifying the final vehicle state in Table 3 is compared to the Dilated Residual Network (DRN) [31]. DRN is a classification model derived from ResNet [32] and is a network that replaces the convolutional layers with dilated convolutional layers. Proposed model utilizes dilated and deformable convolutional layers to extract features from a wide area, and therefore, its network is compared with DRN, which is composed of dilated convolutional layers. DRN has four types, A, B, C, and D,

Table 2. The mAP performance of YOLOv5 on drone train and test dataset.

Class	Images	Instance	mAP@50	mAP@50:95
all_train	9,776	309,470	95.75	83.8
car_vehicle	9,776	277,263	97.2	86.0
truck_vehicle	9,776	32,207	94.3	81.6
all_test	2,200	85,398	91.8	80.3
car_vehicle	2,200	78,765	96.1	85.3
truck_vehicle	2,200	6,633	87.5	75.4

with additional dilated blocks and skip-connections. In this paper, types C and D are used, and type D is a simplified version of type C.

Compared to DRN_D_22, the first proposed model shows a 16.9% difference in accuracy results, but it reduces the number of parameters by 92.2%. In addition, the proposed model presents the results of applying mosaic and color data augmentation. When both augmentations are applied, it shows the best performance among the results presented, with an accuracy of 83.41%. Furthermore, compared to the DRN_C_42 model, it achieves a 1.63% higher accuracy and saves 96% of the parameters.

Table 4 presents the accuracy performance of the proposed model according to the soft label values. The highest accuracy performance of 83.41% is achieved when α is set to 0.7. As α gradually decreases, the performance decreases as well. This is because the soft label values differ from the original hard label values, leading to differences in learning performance.

Table 3. Comparison result with dilated residual networks (DRN) and vehicle state network (VSNet) for data augmentation.

Method	Data augmentation		#para	Acc(%)
	Mosaic	Color		
DRN_C_26	-	-	21,126,584	89.62
DRN_C_42	-	-	31,234,744	81.78
DRN_D_22	-	-	16,393,752	87.69
DRN_D_38	-	-	26,501,912	86.49
DRN_D_54	-	-	35,809,176	89.26
Proposed	-	-	1,273,504	72.85
Proposed	-	O	1,273,504	79.73
Proposed	O	-	1,273,504	81.25
Proposed	O	O	1,273,504	83.41

Table 4. According to Soft label, α , accuracy of vehicle state classification.

Method	Soft label = α	Acc (%)
Proposed	0.3	79.62
Proposed	0.4	81.73
Proposed	0.5	83.12
Proposed	0.6	82.57
Proposed	0.7	83.41

5 Conclusion

This study applies mosaic augmentation and soft-label assignment techniques to classify vehicle states using drone images. Mosaic augmentation combines existing images to create a new image, increasing the amount of data and improving generalization for a limited dataset. Additionally, soft-label assignment is used to generate labels for the mosaic images in vehicle state classification. These two techniques contribute to smooth training and enhance the accuracy performance of the proposed classification model.

Acknowledgment. This work was supported by the National Research Foundation of Korea(NRF) grant funded by the government(MSIT).(No.2020R1A2C200897212).

References

1. Bochkovskiy, A., Wang, C.-Y., Liao, H.-Y.M.: YOLOv4: optimal speed and accuracy of object detection. ArXiv, abs/2004.10934 (2020)
2. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2818–2826 (2015)
3. Cordts, M.: The cityscapes dataset for semantic urban scene understanding. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 3213–3223 (2016)
4. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The KITTI dataset. *Int. J. Rob. Res.* **32**, 1231–1237 (2013)
5. Chen, L.-C., Wang, H., Qiao, S.: Scaling wide residual networks for panoptic segmentation. ArXiv, abs/2011.11675 (2020)
6. Chen, Z.: Vision transformer adapter for dense predictions. ArXiv, abs/2205.08534 (2022)
7. Xu, J., Xiong, Z., Bhattacharyya, S.: PIDNet: a real-time semantic segmentation network inspired from PID controller. ArXiv, abs/2206.02066 (2022)
8. Li, S., Yan, Z., Li, H., Cheng, K.-T.: Exploring intermediate representation for monocular vehicle pose estimation. In: 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1873–1883 (2020)
9. Zhang, Y., Zhang, Q., Zhu, Z., Hou, J., Yuan, Y.: GLENet: boosting 3D object detectors with generative label uncertainty estimation. ArXiv, abs/2207.02466 (2022)

10. Hong, Y., Dai, H., Ding, Y.: Cross-modality knowledge distillation network for monocular 3D object detection. *ArXiv*, abs/2211.07171 (2022)
11. Robicquet, A., Sadeghian, A., Alahi, A., Savarese, S.: Learning social etiquette: human trajectory understanding in crowded scenes. In: Leibe, B., Matas, J., Sebe, N., Welling, M. (eds.) *ECCV 2016*. LNCS, vol. 9912, pp. 549–565. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46484-8_33
12. Zhu, P.: Detection and tracking meet drones challenge. *IEEE Trans. Pattern Anal. Mach. Intell.*, 1 (2020)
13. Cao, Y.: VisDrone-DET2021: the vision meets drone object detection challenge results. In: *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pp. 2847–2854 (2021)
14. Du, D., et al.: VisDrone-CC2021: the vision meets drone crowd counting challenge results, pp. 2830–2838 (2021)
15. Fan, H.: VisDrone-MOT2021: the vision meets drone multiple object tracking challenge results. In: *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pp. 2839–2846 (2021)
16. Wang, J., Xu, C., Yang, W., Yu, L.: A normalized Gaussian Wasserstein distance for tiny object detection. *ArXiv*, abs/2110.13389 (2021)
17. Lee, Y., Tang, Q., Choi, J.-W., Jo, K.: Low computational vehicle re-identification for unlabeled drone flight images. In: *IECON 2022–48th Annual Conference of the IEEE Industrial Electronics Society*, pp. 1–6 (2022)
18. Lee, Y., Tang, Q., Choi, J.-W., Jo, K.: Low computational vehicle lane changing prediction using drone traffic dataset. In: *2022 International Workshop on Intelligent Systems (IWIS)*, pp. 1–4 (2022)
19. Krajewski, R., Bock, J., Kloeker, L., Eckstein, L.: The highD dataset: a drone dataset of naturalistic vehicle trajectories on German highways for validation of highly automated driving systems. In: *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2118–2125 (2018)
20. Bock, J., Krajewski, R., Moers, T., Runde, S., Vater, L., Eckstein, L.: The inD dataset: a drone dataset of naturalistic road user trajectories at German intersections. In: *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1929–1934 (2019)
21. Krajewski, R., Moers, T., Bock, J., Vater, L., Eckstein, L.: The round dataset: a drone dataset of road user trajectories at roundabouts in Germany. In: *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6 (2020)
22. Jocher, G.R.: ultralytics/yolov5: v5.0 - YOLOv5-P6 1280 models, AWS, supervise.ly and YouTube integrations (2021)
23. Lin, T.-Y., et al.: Microsoft COCO: common objects in context. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014*. LNCS, vol. 8693, pp. 740–755. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-10602-1_48
24. Lee, Y., Kim, S., Choi, J., Jo, K.: Vehicle state classification from drone image. In: *IEEE International Conference on Industrial Technology (ICIT)*, pp. 1–5 (2023)
25. Ioffe, S., Szegedy, C.: Batch normalization: accelerating deep network training by reducing internal covariate shift. In: *International Conference on Machine Learning* (2015)
26. Hendrycks, D., Gimpel, K.: Gaussian error linear units (GELUs). [arXiv: Learning](https://arxiv.org/abs/1606.02256) (2016)
27. Dai, J.: Deformable convolutional networks. In: *2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 764–773 (2017)

28. Woo, S., Park, J., Lee, J.-Y., Kweon, I.S.: CBAM: convolutional block attention module. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds.) ECCV 2018. LNCS, vol. 11211, pp. 3–19. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-01234-2_1
29. Lin, T.-Y., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal loss for dense object detection. *IEEE Trans. Pattern Anal. Mach. Intell.* **42**, 318–327 (2017)
30. Kingma, D.P., Ba, J.: Adam: a method for stochastic optimization. *CoRR*, abs/1412.6980 (2014)
31. Yu, F., Koltun, V., Funkhouser, T.: Dilated residual networks. In: *Computer Vision and Pattern Recognition (CVPR)* (2017)
32. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778 (2016)