

Vehicle Perception from Satellite

Bin Zhao, *Member, IEEE*, Pengfei Han, Xuelong Li, *Fellow, IEEE*

Abstract—Satellites are capable of capturing high-resolution videos. It makes vehicle perception from satellite become possible. Compared to street surveillance, drive recorder or other equipments, satellite videos provide a much broader city-scale view, so that the global dynamic scene of the traffic are captured and displayed. Traffic monitoring from satellite is a new task with great potential applications, including traffic jams prediction, path planning, vehicle dispatching, *etc.*. Practically, limited by the resolution and view, the captured vehicles are very tiny (a few pixels) and move slowly. Worse still, these satellites are in Low Earth Orbit (LEO) to capture such high-resolution videos, so the background is also moving. Under this circumstance, traffic monitoring from the satellite view is an extremely challenging task. To attract more researchers into this field, we build a large-scale benchmark for traffic monitoring from satellite. It supports several tasks, including tiny object detection, counting and density estimation. The dataset is constructed based on 12 satellite videos and 14 synthetic videos recorded from GTA-V. They are separated into 408 video clips, which contain 7,336 real satellite images and 1,960 synthetic images. 128,801 vehicles are annotated totally, and the number of vehicles in each image varies from 0 to 101. Several classic and state-of-the-art approaches in traditional computer vision are evaluated on the datasets, so as to compare the performance of different approaches, analyze the challenges in this task, and discuss the future prospects. The dataset is available at: <https://github.com/Chenxi1510/Vehicle-Perception-from-Satellite-Videos>.

Index Terms—remote sensing, tiny object detection, vehicle counting, density estimation.



1 INTRODUCTION

Recently, with the significant progress of aerospace technology, the commercial satellites are able to capture the Very High Resolution (VHR) videos [1], [2]. For example, the Jilin-1 satellites can observe the earth with the spatial resolution of 0.72m [3]. SkySat-1 satellites provide the VHR videos with the resolution around 1m [4]. These satellites can monitor the ground dynamically in city-scale, where the vehicles can be seen clearly [5]. VHR videos from satellites provide a new perspective for traffic monitoring.

Traffic monitoring from satellite is quite different from those street surveillance cameras on the ground. It has a variety of advantages:

1) **Geographical Boundlessness:** Satellite transcends terrestrial geographical constraints, encompassing urban, rural, and remote terrains. It orbits dynamically, affording a panoramic vantage point for comprehensive vehicular oversight. In contrast, surveillance cameras are constrained to urban centers and pivotal transport hubs. The installation of surveillance cameras worldwide is prohibitively costly due to extreme geographical conditions. Yet, such coverage is pivotal for traffic monitoring in remote areas and holistic network analysis.

2) **Real-time and Periodicity:** Satellites furnish nearly real-time imagery and data, facilitating immediate traffic monitoring and emergency response. Furthermore, routine satellite missions yield long-term time-series data for traffic trend and pattern analysis. They swiftly monitor traffic conditions during natural disasters or emergencies, providing crucial information for rescue operations.

3) **Wide-area Surveillance:** Satellite traffic monitoring spans vast geographical regions, from cities to broader territories. This has immense value for urban planning, traffic management, and natural disaster monitoring, offering a global perspective. In contrast to ground-based surveillance cameras that can only observe traffic in intersections, streets, or roundabouts, satellites can monitor sprawling areas spanning square kilometers. Their expansive field of view enables the surveillance of city-scale



Fig. 1: Examples of (a) satellite images captured over Dubai International Airport, and (b) the 1,214 vehicles in the scene highlighted by yellow points.

traffic, providing a fresh perspective for traffic control, analysis, and planning.

4) **Monitoring Diverse Transport Modes:** Satellites can monitor not only road traffic but also railways, aviation, and maritime transport, fostering multimodal traffic research. Additionally, satellites can complement surveillance cameras, enabling layered traffic monitoring in cities. Surveillance cameras offer localized and detailed insights, while satellites provide global traffic monitoring. They can be integrated to monitor the traffic from the ground view and sky view hierarchically.

Apart from the above advantages, traffic monitoring from satellite is also a challenging task:

1) The vehicles in satellite videos are tiny. Limited by the resolution, the vehicles only contain a few pixels (10 pixels or so) and lack of appearance information. Some of them cannot even be recognized unless they are moving, as depicted in Fig. 1. The detection and counting of such tiny vehicles are very difficult. Traditional object detection and counting approaches cannot deal with them effectively and efficiently.

2) The movement in satellite videos is very complicated. As depicted in Fig. 2 (a) and (b), the background and vehicles are both

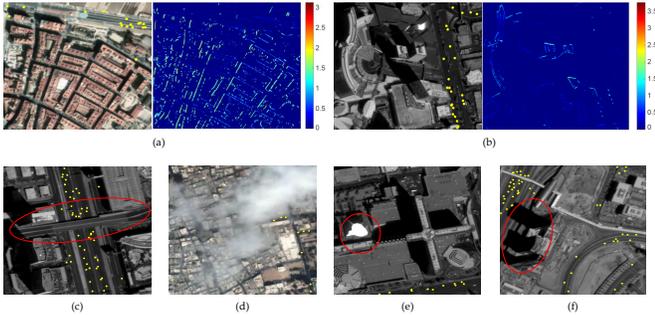


Fig. 2: Challenges in vehicle perception from satellite. The optical flow maps in (a) and (b) indicate the movements are complex and uneven. (c), (d), (e) and (f) display the noise in satellite videos. They are shelter, clouds, specularity and shadow from left to right.

moving. Practically, the movement of background is correlated to the view of satellite and the height of buildings, since the video only captures the 2D projection of the 3D movement. As a result, it is hard to separate the moving vehicles from the background. Moreover, the gradual movement of the satellite engenders localized misalignment and dynamic intensity variations in some stationary background objects. These changes serve no purpose in discerning the motion of objects and may potentially give rise to motion artifacts.

3) The satellite videos are full of noise. As evidenced in Fig. 2 (c)-(f). The frames captured are subject to an array of perturbations, including obstructions from edifices, clouds in the sky, shadows cast by sunlight, and mirrored reflections on ground-level glass surfaces. Under such circumstances, the local contrast between the background and targets diminishes significantly. Owing to the complexity of the background and the presence of image noise, targets occasionally amalgamate with the chaotic background, resulting in localized obscurity. This noise exerts substantial interference in the detection of vehicles.

4) The LEO satellites are incapable of perpetually hovering over a certain urban domain. In such instances, satellites are unable to sustain prolonged and continuous traffic surveillance. Additionally, when the satellite platform is fixated on a specific area of interest, the acquired satellite video evinces localized positional discrepancies and alterations in local intensity, attributable to immobile objects. These objects may be erroneously discerned as mobile entities, thereby exacerbating the incidence of false positives in the data.

1.1 Motivation

Traffic monitoring from satellite is still in the preliminary stage. The lack of released satellite videos and large-scale annotated datasets are the key factors to limit the development. To the best of our knowledge, only two videos are released with partially annotations [1], [6]. The lack of data leads to the following dilemmas in the research:

1) Most of existing approaches just evaluate their performance on one or two videos. Such little test samples are not enough to verify the effectiveness. Worse still, some of the annotations or videos are not released publicly, which causes interference to the fair comparison of different approaches.

2) Deep learning approaches are not applicable to this field. Most existing approaches are developed based on traditional detectors, since the annotated samples are not enough for training.

However, deep learning has become the mainstream and surpasses traditional approaches in most computer vision tasks. The performance of this field is impeded by the lack of annotated data.

1.2 Overview

Practically, the annotations of vehicles in satellite videos are difficult, since the frames are in city-scale, the movements of background are uneven, the sunlights are varying, and the vehicles are very tiny and cannot even be recognized by human eyes. In this paper, we make efforts to overcome these difficulties, and construct a large-scale dataset for traffic monitoring from satellite, which is named as TMS. TMS is composed of 408 videos collected from 12 real satellite videos and 14 synthetic videos of GTA-V, where the synthetic videos are utilized to make up for the lack of available satellite videos. In each video, the coordinates of vehicles are annotated at 1 frame per second. The numbers of vehicles in each frame vary from 0 to 101, and 128,801 vehicles are annotated totally. With the help of the TMS dataset, three tasks are developed for traffic monitoring from satellite, *i.e.*, Tiny Object Detection (TOD), VEHICLE Counting (VEC) and Traffic Density Estimation (TDE). Numbers of classic and state-of-the-art approaches are evaluated, including both traditional and deep learning approaches, where the challenges of each task are analyzed and the insights for researchers are presented.

1.3 Contributions

Overall, the main contributions of this paper can be summarized as follows:

1) The largest satellite video dataset is constructed for traffic monitoring. It can promote the research in this field by attracting deep learning approaches and provide an evaluation platform for different approaches.

2) The synthetic videos from GTA-V are integrated with real videos to further augment the scale of real dataset. They are annotated automatically, and provide a new perspective to relieve the problem of lacking real data.

3) Three tasks are performed on this benchmark, including Tiny Object Detection (TOD), VEHICLE Counting (VEC), and Traffic Density Estimation (TDE), so as to promote the development of traffic monitoring from satellite.

2 RELATED WORKS

In the following subsections, the related tasks in remote sensing, computer vision, and the recently proposed approaches of traffic monitoring from satellite are reviewed.

2.1 Object Detection in Remote Sensing

Object detection is a long standing task in remote sensing, which is also the basis of traffic monitoring from satellite. Numbers of datasets are constructed, *e.g.*, TAS [7], SZTAKI-INRIA [8], NWPU VHR-10 [9], HRSC2016 [10], DOTA [11], *etc.*. The images are mainly collected from satellite, aerial plane, and other platforms, such as Google Earth, Tianditu and Quickbird. The annotated objects are in multiple categories. Vehicle is the most popular category for object detection, which shows its importance in remote sensing [12]. Earlier object detection approaches are developed based on template matching, geometry modeling, context knowledge, and low-level feature extraction [13]. However, object detection in remote sensing is quite complex with the

factors of noise, size, lights and background [14]. These traditional approaches are not robust enough to be generalized to different situations. Recently, deep learning is employed in this task under the support of large-scale datasets. Most of them are modified from mainstream object detection approaches in natural scene images, *e.g.*, Faster RCNN [15], SSD [16], YOLO [17]. By taking advantages of the non-linear learning ability, deep learning approaches surpass traditional approaches, and boost the performance significantly.

2.2 Relevant Tasks in Computer Vision

Object detection is a classic computer vision task for natural scene images. Traditional approaches follow a pipeline of region selection (*e.g.*, superpixels [18], sliding window [19], selective search [20]), feature extraction (*e.g.*, SIFT [21], HOG [22]) and classifier (*e.g.*, SVM [23], Adaboost [24]). However, traditional approaches and hand-crafted features are not capable enough to generalize to the variance of size, shape, occlusion and noise in object detection. Till now, deep learning approaches are taking the leading position, including RCNN series [15], [25], YOLO series [26], [27], [28], *etc.* They are mostly developed based on challenges, *i.e.*, PASCAL VOC¹, ILSVRC² and MS-COCO Detection³. Object detection in natural images can inspire vehicle detection in satellite videos. Object counting and density estimation are important for crowd analysis as well as traffic congestion prediction. Apart from adopting object detectors directly, CNN-based regression models are widely used in this task. Most of them are developed based on Fully Convolutional Network (FCN) [29] with the instance-level or image-level annotations. TRANCOS [30] and VisDrone2019 Vehicle [31] are two datasets specially designed for vehicle counting and density estimation, which can benefit traffic monitoring from satellite.

2.3 Traffic Monitoring from Satellite

With the release of several satellite videos from Jilin-1 and SkySat-1, traffic monitoring from satellite draws increasing attention [1], [6]. Most works are presented in recent years. In [32], the vehicles are detected by background subtraction, where mathematic morphology and statistical analysis are utilized to estimate background of each frame. The performance is evaluated on one satellite video from SkySat-1. The low-rank matrix decomposition is modified in [4], in order to model background and foreground with the regularization of low-rank and sparsity. Furthermore, matrix decomposition assisted with moving confidence is developed in [1], which can promote the motion of vehicles meanwhile suppress that of the background. The performance is evaluated on two satellite videos and two surveillance videos. Overall, traffic monitoring from satellite is still in its early stages, and requires large-scale datasets and benchmarks to advance the development.

3 THE TMS DATASET

3.1 Data Collection and Preprocessing

We tried our best to collect the satellite videos that can be used for traffic monitoring from satellite. The real part of TMS is composed of 12 full satellite videos. They are captured by the

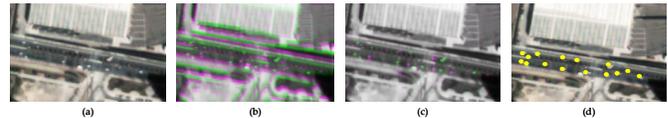


Fig. 3: The annotation process of each real image in TMS. (a) is the original image. (b) shows the difference between consecutive frames. (c) is the difference of registered consecutive frames. (d) displays the final annotation of vehicles.

non-stationary satellite platform, *i.e.*, Jilin-1 and SkySat-1, with the ground sample distance of 1m or so. In this case, the vehicles only contain 5-20 pixels. The videos are recorded over the sky of Boston, Dubai, Valencia, Jeddah, LasVegas, Hong Kong, Aleppo, Bangkok, and Tokyo. The captured scenes include city street, airport, suburbs, and port. Five of the videos are freely provided by the official website of Jilin-1⁴, and the others are crawled from YouTube⁵, since the officially released versions are not available.

The synthetic part of TMS contains 14 videos obtained from the game of GTA-V. GTA-V builds a fictional city with Los Angeles as the prototype, namely, Los Santos. The map in the game covers 252 square kilometers. The scene rendering, lights, shadow, weather effects and other conditions are quite similar to those in the real world, so that the players can immerse themselves into it. Furthermore, GTA-V is allowed to be developed by players for non-commercial use, such as academic use. Different from capturing videos from the egocentric view of the player, the videos are recorded from the satellite view of GTA-V in this paper. These videos display different city streets or suburbs in the map of the game. They are recorded by the screen recording software of Windows 10, *i.e.*, Xbox Game Bar.

The real and synthetic satellite videos are combined together to form the TMS dataset. Each video is with the resolution of 1080x1920 or 3072x4096, which covers an area of several square kilometers. For simplicity, these full satellite videos are segmented spatially into subareas. In this case, those videos with the resolution of 1080x1920 are segmented evenly into 540x480 subareas by 2x4 grid, and those videos with the resolution of 3072x4096 are segmented evenly into 512x512 subareas by 6x8 grid. Totally, 408 videos are obtained by segmenting full videos into subareas, including 296 real videos and 112 synthetic videos. Besides, to simplify the annotation, each video is sampled by 1 fps. Finally, TMS contains 9,296 images. The number of real images and synthetic images are 7,336 and 1,960, respectively.

The annotation of satellite videos is quite arduous. It is because the vehicles are tiny and lack of appearance features, so that they can hardly be recognized from background. To address this problem, a motion based annotation method is conducted in our work. Specifically, the motion of vehicles is larger than the surrounding background. In this case, the vehicles can be distinguished from background by comparing the difference between two consecutive frames. To amplify the difference, two frames with the temporal interval of one second are compared, so that the difference can be recognized by human eyes. As depicted in Fig. 3, the differences are highlighted by colors. The green color means the pixel value of current frame at this region is higher than the next frame, which indicates the location of vehicles. In contrast, the red color indicates the vehicle locations of next frames. However, the

1. <http://host.robots.ox.ac.uk/pascal/VOC/index.html>

2. <https://image-net.org/challenges/LSVRC/2017/>

3. <https://cocodataset.org/#home>

4. <https://mall.charminglobe.com/Sampledata>

5. <https://www.youtube.com/>

TABLE 1: Statistics of reported datasets of traffic monitoring from satellite.

Datasets	#Videos	#Resolution	#Images	#Vehicles			Tasks			Availability
				Min	Max	Total	TOD	VEC	TDE	
SHDV [33]	1	400*400	700	–	–	–	✓	✗	✗	No
Valencia [6]	3	500*500	168	7	41	3,211	✓	✗	✗	Yes
Las Vagas [4]	2	400*400, 600*400	1,400	27	86	80,047	✓	✗	✗	Yes
Jilin-1 [1]	3	(400~700)*(400~600)	900	–	–	–	✓	✗	✗	No
SkySat [1]	6	(400~600)*(400~600)	3,500	–	–	–	✓	✗	✓	Partially
TMS	408	512*512, 540*480	9,296	0	101	128,801	✓	✓	✓	Yes

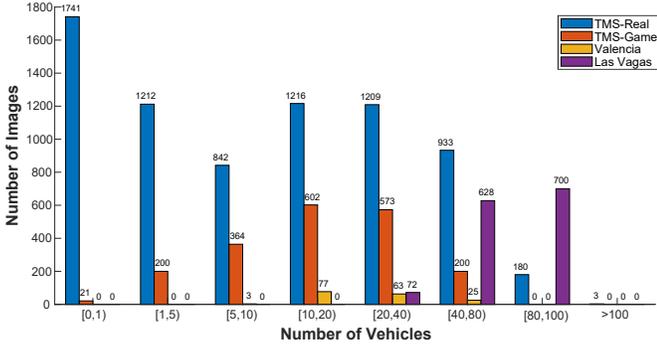


Fig. 4: Image distribution of different datasets over the number of vehicles. Note that TMS-Real and TMS-Game represent the real and synthetic parts of TMS, respectively.

motions from vehicles and background are intertwined. It makes the annotation in a complete daze, as shown in Fig. 3 (b).

To further reduce the interference caused by background motion, the translation of background is removed by the intensity-based image registration method⁶. In this case, the annotator can localize the vehicles effortlessly. Each vehicle is annotated with a point which indicates the location. Practically, the annotation tool is developed based on Matlab R2019b.

In synthetic videos, the locations of vehicles are obtained automatically by transforming the map coordinate into the screen coordinate. The annotation tool is designed based on the game plugin developer, Script Hook V [34]. It can save much labor force in the annotation process. Besides, it can generate lots of automatically annotated synthetic videos, in order to augment the training data and boost the performance on real satellite videos. It provides an efficient way to remedy the lack of real data in VHR satellite videos.

We want to emphasize that, different from the previously released Valencia dataset [6], the IDs of vehicles in different frames are not provided in the proposed TMS dataset⁷. It is for the following reasons: 1) Valencia is a partially annotated dataset, where the annotated regions are selected manually. In contrast, TMS is a fully annotated dataset. The vehicles are very dense and lack of appearance features. It can hardly be identified by annotators. 2) The VHR videos are captured by satellites in low earth orbit. The view of satellite is non-stationary, so that it can only gaze the city at a quite short time, which makes it infeasible to track vehicles effectively. In this case, the vehicle tracking task is not developed in this paper.

Overall, in our work, 14 annotators are employed in the

6. <https://www.mathworks.com/help/images/ref/imregister.html#description>

7. Practically, our plugin is able to identify the vehicles in different frames.

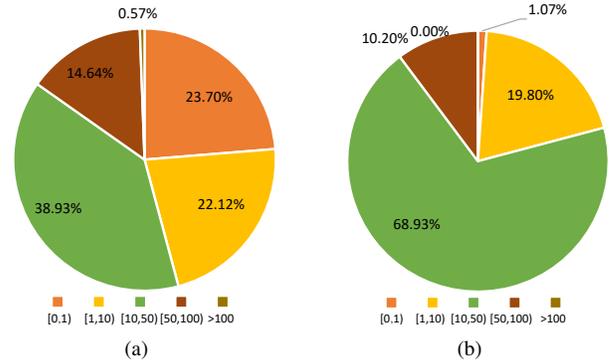


Fig. 5: Vehicle density distribution of TMS-real (a) and TMS-game (b).

annotation process, and more than 500 working hours are spent. Furthermore, each annotator also plays the role of checker, in order to guarantee the correctness of annotations for each image.

3.2 Data Statistics

The TMS dataset is constructed from 26 real and synthetic satellite videos, where 408 video clips and 9,296 images are obtained. In each image, the locations of vehicles are annotated, and the number of vehicles varies from 0 to 101. In total, 128,801 vehicles are annotated.

Table 1 presents the statics of existing datasets. At present, TMS is the largest dataset for traffic monitoring from satellite. It is also the only one publicly available large-scale dataset in this task. Compared to existing datasets, it has the following advantages:

1) The real part of TMS is created based on 12 fully annotated satellite videos, each covers an area of several square kilometers, and 7,336 annotated images are obtained. However, existing datasets only contain one or two partially annotated videos. With the large-scale TMS dataset, the deep learning approaches are applicable to the tasks of traffic monitoring from satellite. The evaluation can also be conducted effectively.

2) TMS is the first to utilize synthetic videos to augment the dataset. It can address the problem of lacking of publicly available satellite videos. Synthetic videos can be annotated automatically, which can reduce the labor costs significantly. By utilizing synthetic videos, TMS pioneers researchers a new perspective to develop approaches and boost the performance.

3) In TMS, the videos are captured over different cities all around the world as well as the virtual city. The scenes include downtown, suburbs, harbor, airport, *etc.*. The vehicles in each image vary in a large range from 0 to 101. However, existing datasets only contain one or two scenes with dense traffic, which

cannot meet the demand of traffic monitoring in different places and situations. In this case, TMS can improve the generality of developed approaches.

Furthermore, Fig. 4 and 5 plot the distribution of samples over the number of vehicles and traffic density. It can be observed that TMS holds the largest range over the number of vehicles among different datasets, *i.e.*, from 0 to 100+. Besides, only TMS contains negative samples in which the number of vehicles is zero. The other two datasets, Valencia and Las Vegas, lack of null and sparse traffic scenes, since most of the samples are with the vehicles over 10. Furthermore, we can see from Fig. 5 that both the real part and synthetic part of TMS cover the null, sparse, normal and dense traffic scenes. It shows the superiority of TMS compared with existing publicly available datasets. In conclusion, TMS is more applicable to the real scenarios of traffic monitoring from satellite.

3.3 Applicable Tasks

To monitor traffic from satellite effectively, three tasks are developed based on the TMS dataset. They are Tiny Object Detection (TOD), Vehicle Counting (VEC) and Traffic Density Estimation (TDE). For the TMS dataset, we randomly allocated 50% as the training set, 25% as the validation set, and the rest for testing, adhering to this criterion across all tasks. Note that the dataset is allocated in video clip-wise (*i.e.*, 408 clips totally) while not image wise, since two temporally adjacent images are similar and may result in training information leakage.

The target of TOD is to localize the vehicles in each image. It can be conducted by two strategies. One is the single-image detection strategy, where the objects are detected by applying detectors to single image. The other is motion-based detection strategy. The detectors operate on the video. The motion among consecutive frames is utilized to distinguish objects from background. Following existing detection protocols, the annotated points are modified to bounding boxes with the height and width of six pixels. The size of bounding boxes is large enough to cover most vehicles. In the evaluation, Precision, Recall and F-score are utilized as the metrics.

The task of VEC is to count the number of vehicles in the screen. It is a basic factor to estimate the traffic congestion. Different from TOD, VEC focuses on the global traffic situation, which is not necessary to localize each vehicle in the image. In this paper, the Mean Average Error (MAE) and Mean Square Error (MSE) are adopted as the evaluate metric, which can measure the difference of the estimated number and annotated number.

TDE aims to estimate the traffic density map in each image. It provides a vivid visualization of traffic congestion. In practice, to generate the ground truth of density maps, the annotated location marks are blurred by the Gaussian kernel, where the kernel size is fixed as 29, and $\sigma = 4$. Two popular metrics are adopted in TDE, *i.e.*, Structural Similarity in Image (SSIM) and Peak Signal-to-Noise Ratio (PSNR). They measure the performance of density map estimation.

4 EXPERIMENTS

4.1 Experiments on Tiny Object Detection

In the TOD task, three kinds of approaches are evaluated, *i.e.*, background subtraction, deep object detector, deep object localizer. They are described in the following subsections.

4.1.1 Background Subtraction

R-PCA: Robust Principle Component Analysis [35]. It separates the background and foreground into a low-rank matrix and a sparse matrix, respectively. It is optimized jointly by the principle component pursuit and fast low-rank approximation. The foreground pixels are obtained by morphological segmentation.

GMM: Gaussian Mixture Model [36]. It is developed with the assumption that the each pixel follows the Gaussian distribution temporally. It models the background with several Gaussian distributions. One pixel is identified as the background point once it belongs to one of distributions, and the background is updated iteratively.

ViBe: Visual Background extractor [37]. It models each background pixel with several neighbors. In this case, the background can be initialized in single frame. With the assumption that stochastic models can simulate the uncertainty of pixel change, the background is updated with one random neighbor if one point is identified as the background.

TFD: Three Frame Difference [38]. It subtracts the background by detecting the changes among frame sequence. Three consecutive frames are utilized to compute the changes pairwise, which can effectively reduce the interference caused by occlusions and noise.

LSD: Low-rank and Structured sparse Decomposition [39]. It decomposes the foreground and background with the regularization of low-rank norm and structured sparse norm.

DECOLOR: DETecting Contiguous Outliers in the LOW-rank Representation [40]. It represents the video frames as a low-rank matrix with the assumption that the background are linearly correlated, and objects are detected as outliers.

4.1.2 Deep Object Detector

Faster RCNN: Faster Region-CNN [15]. It is a popular anchor-based deep object detector. It firstly utilizes the Region Proposal Network (RPN) to generate region proposals, and then apply the Region-CNN to regress the coordinate of detected objects.

SSD: Single Shot multibox Detector [16]. It is a one-stage object detector that localizes and recognizes objects jointly. It is much efficient than Faster RCNN. In this paper, VGG-16 is adopted as the backbone of SSD.

YOLO series: The fourth, fifth, sixth, seventh, and eighth version of You Only Look Once [28]. YOLO is an anchor-free deep object detector. It segments the whole image into grid cells, and regresses the localization directly. YOLO v4 [28], YOLO v5⁸, YOLO v6 [41], YOLO v7 [42], and YOLO v8 [43] are quite similar. They both utilize techniques, such as cropping, rotation, flipping, mosaic, *etc.*, to augment the training data.

Tood: Task-aligned one-stage object detection [45]. It is a one-stage object detector [45], which explicitly aligns object classification and localization in a learning-based manner. It utilizes ResNet as the backbone and follows an overall pipeline of backbone-FPN-head.

Dino: DETR with Improved deNoising anchor boxes [44]. It is a robust end-to-end approach for improved denoising anchor boxes in object detection, and Transformer is adopted as the backbone.

4.1.3 Deep Object Localizer

RAZNet: Recurrent Attentive Zooming Network [47]. It is originally proposed for crowd localization. In this paper, it is utilized to localize the vehicles just with the point annotations. It is good

8. <https://github.com/ultralytics/yolov5>

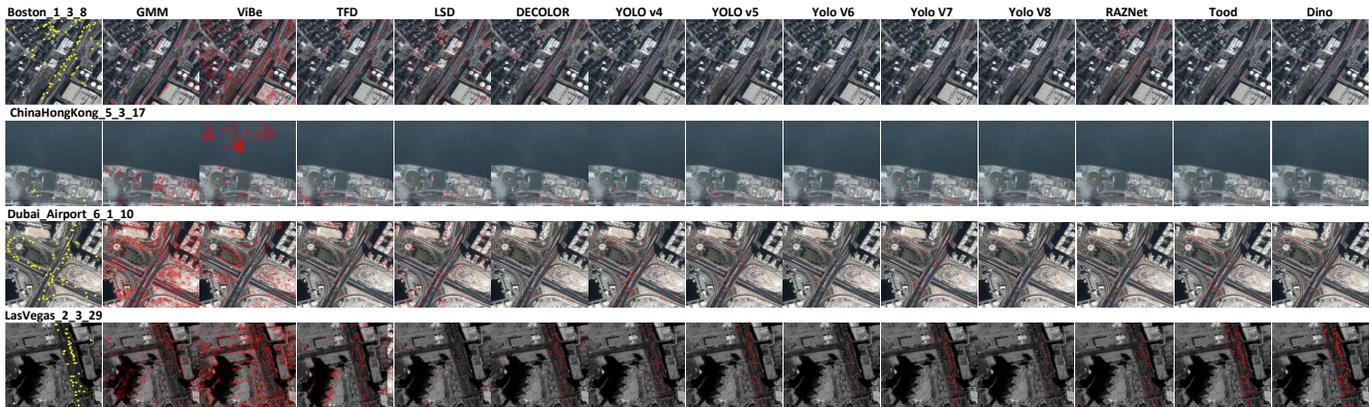


Fig. 6: Visualization exemplar results of different approaches on the TOD task across the TMS dataset. The first column denotes the annotated detection results, and the other columns represent the final detection results of different methods.

TABLE 2: TABLE II TOD RESULTS OF DIFFERENT APPROACHES ON TMS. THE BEST RESULTS ARE MARKED IN BOLD.

Setting	Real			Game			Augmented		
	<i>Precision</i> ↑	<i>Recall</i> ↑	<i>F_{score}</i> ↑	<i>Precision</i> ↑	<i>Recall</i> ↑	<i>F_{score}</i> ↑	<i>Precision</i> ↑	<i>Recall</i> ↑	<i>F_{score}</i> ↑
R-PCA [35]	6.67	0.59	0.97	1.43	4.64	2.03	---	---	---
GMM [36]	23.60	3.07	5.19	15.09	4.18	6.15	---	---	---
ViBe [37]	54.63	9.51	14.92	30.68	17.56	21.49	---	---	---
TFD [38]	69.33	9.28	15.63	7.96	2.69	3.89	---	---	---
LSD [39]	46.15	13.61	19.61	38.00	29.76	32.98	---	---	---
DECOLOR [40]	52.48	56.65	52.67	32.21	46.01	37.23	---	---	---
Faster RCNN [15]	2.56	3.92	3.00	1.68	7.13	3.00	1.91	3.39	2.00
SSD [16]	3.20	3.55	1.00	0.25	3.41	0.63	0.32	3.58	1.00
YOLO v4 [28]	35.99	14.75	22.00	39.41	15.71	22.00	45.21	23.52	31.00
YOLO v5	39.8	27.1	32.2	24.2	16.3	26.1	40.6	28.5	33.5
YOLO v6 [41]	35.11	31.3	36.43	37.63	17.85	31.06	43.5	35.2	37.6
YOLO v7 [42]	28.03	31.1	29.45	22.4	17.2	19.47	31.6	23.5	33.2
YOLO v8 [43]	35.61	32.71	39.81	40.13	19.63	25.82	47.1	28.4	35.9
Dino [44]	5.4	45.9	9.7	6.3	40.9	10.9	6.92	49.1	12.3
Tood [45]	14.0	40.3	20.7	8.8	33.1	13.9	19.3	46.5	22.8
SCAL_Net [46]	0.34	6.92	2.17	0.33	0.19	0.24	0.37	5.47	2.16
RAZNet [47]	42.77	53.24	47.53	28.05	59.04	38.03	51.31	62.12	56.78

at detecting tiny objects by operating recursively on small image regions and zooming them into high-resolution.

SCAL_Net: A Simple yet effective Counting And Localization Network [46]. It proposes a joint framework for vehicle counting and localization, which tackles them as a pixel-wise dense prediction problem.

4.1.4 Results and Discussions

Table 2 presents the results of different approaches on the task of tiny object detection. It can be observed that TOD is really a challenging task, since most of the classic and popular approaches perform worse than they are on traditional object detection tasks. For background subtraction approaches, R-PCA gets really poor performance. It is mainly because the foreground in R-PCA is obtained by morphological segmentation, while these vehicles are so tiny that their morphological features can hardly be recognized. Worse still, other background modeling approaches, including GMM, ViBe and LSD, also suffer from the extremely unbalanced distribution of foreground and background pixels. TFD distinguishes vehicles from background by the changes among three consecutive frames. However, as aforementioned, the motion in satellite videos are quite complex. DECOLOR takes vehicles as outliers with the assumption that the background are linearly

correlated in the frame sequence, which performs the best in background subtraction approaches.

Deep learning based object detection has achieved great success in the computer vision field. However, TOD is an extremely challenging task, since the vehicles in satellite videos are too tiny to be captured by the detector. It can be observed that Faster RCNN and SSD get poor performance, even worse than traditional detectors. It is mainly because that the scales of vehicles only take a minor part of the whole image. The appearance of vehicles is vanished from feature maps after the encoding of several convolutional and pooling layers. To maintain the vehicle features, YOLO v4, v5, v6, v7, and v8 utilize a more powerful and deeper backbone network to extract richer feature representations. In this case, their performance surpasses that of Faster RCNN and SSD, as well as outperforms Dino and Tood, which utilize multi-scale feature interactions.

The results of two deep object localizers are also presented in Table 2. They detect vehicles from point annotations, which is quite different from object targeting approaches. We can see that SCAL_Net performs poor. It is because the vehicles are tiny and much smaller than the instance in the crowd. RAZNet achieves much better results on the real part of TMS. It is mainly benefited from its zooming strategy, which is operated on small regions and zooming into high-resolution recursively.

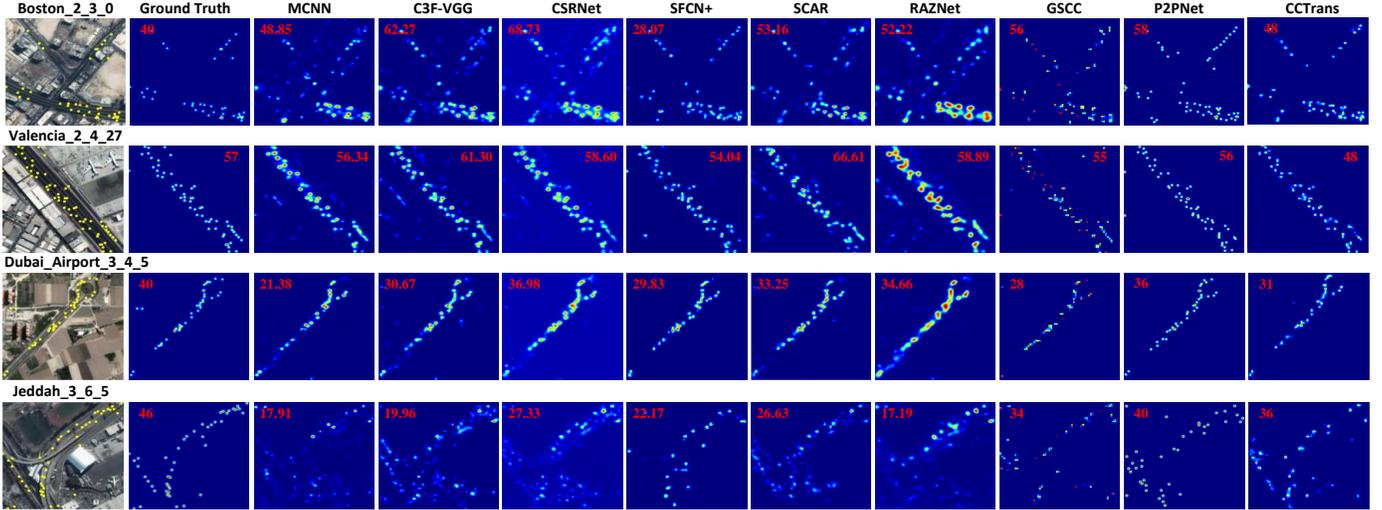


Fig. 7: Visualization of the crowd density maps of different methods on the VEC and TDE tasks across the TMS dataset. The first and second columns denote the annotated density maps, and the other columns represent the final density estimation results of different methods. Note that the heat of each pixel indicates the vehicle density. The red numbers denote the annotated and predicted counting results.

TABLE 3: VEC AND TDE RESULTS OF DIFFERENT APPROACHES ON TMS. THE BEST RESULTS ARE MARKED IN BOLD.

Setting	Real		Game		Augmented	
	MAE/MSE ↓	PSNR/SSIM ↑	MAE/MSE ↓	PSNR/SSIM ↑	MAE/MSE ↓	PSNR/SSIM ↑
MCNN [48]	7.15/113.83	29.02/0.90	7.28/107.37	30.89/0.93	6.92/106.43	29.33/0.90
C3F-VGG [49]	6.99/97.75	29.09/0.90	10.23/189.95	30.37/ 0.97	5.98/83.25	29.31/0.93
CSRNet [50]	6.75/139.18	29.31/0.89	10.46/195.96	31.03/0.94	6.02/77.83	29.36/0.94
SFCN+ [51]	4.75/70.62	28.95/ 0.97	11.28/234.42	30.65/ 0.97	12.34/287.61	28.66/0.79
SCAR [52]	6.70/146.64	29.70/0.96	8.01/125.58	30.58/0.94	5.79/105.49	29.71/0.96
RAZNet [47]	7.17/102.47	30.17/0.94	7.85/130.33	29.79/ 0.97	6.74/96.57	31.28/0.97
SCAL_Net [46]	13.02/354.57	8.74/0.56	10.87/182.47	17.57/0.78	12.97/329.33	9.23/0.67
GSCC [53]	5.19/78.61	28.22/0.87	5.28/97.39	29.55/0.89	4.89/78.05	29.25/0.91
P2PNet [54]	5.96/117.36	29.87/0.91	5.38/105.1	31.42/0.96	5.03/ 74.29	30.76/0.95
CCTrans [55]	5.34/ 84.37	29.57/0.92	4.42/59.62	29.87/0.95	5.96/87.35	29.54/0.92

Fig. 6 presents the results of different approaches on the TOD task. Note that those approaches perform extremely poor in Table 2 are not displayed. We can see that the results of GMM, ViBe and TFD contain lots of false positive samples. Others performs better but far from satisfactory. Overall, from the results in Table 2 and Fig. 6, we can simply conclude that TOD is a challenging task due to the tiny size of vehicles. As a result, most traditional approaches get much worse performance than they are on the typical object detection task. Besides, deep learning based detectors and localizers show less superiority as well. Fortunately, spatial zooming and feature fusion provide some meaningful insights for future researches.

4.2 Experiments on Vehicle Counting and Traffic Density Estimation

Vehicle counting and density estimation are two relevant tasks. They follow similar protocols in the literature. In this case, the experiments of VEC and TDE are conducted and discussed in this subsection together.

4.2.1 Methods for Comparison and Evaluation

MCNN: Multi-Column CNN [48]. By designing kernel in different sizes, a multi-column scheme is conducted, which can adapt to arbitrary crowd density and perspective map of images.

C3F-VGG: A variant of VGG-16 [49]. The first 10 convolutional layers of VGG-16 are adopted for representation learning, and another two regression layers are added to estimate the density map.

CSRNet: Congested Scene Recognition Network [50]. It also takes VGG-16 as the backbone, and designs a dilation module on the top.

SFCN+: Spatially Fully Convolutional Network [51]. It takes ResNet-101 [56] as the backbone, and adds dilation convolution, spatial encoder and regression to predict the density map directly.

SCAR: Spatial-Channel-wise Attention Regression [52]. It operates self-attention on the spatial axis and channel axis of feature map, which can capture the global contextual information.

GSCC. Generalized losS function for Crowd Counting [53]. It learns density map representations through the optimization of the non-equilibrium problem, and introduces a generalized loss function for acquiring density maps used in crowd counting and localization.

P2PNet. Point to Point Network [54]. It presents a purely point-based framework, allowing for a more precise and seamless integration of individual localization with crowd counting.

CCTrans. Crowd Counting with Transformer [55] It leverages a pyramid visual transformer backbone to capture overarching crowd information, followed by the utilization of another pyramid visual transformer backbone to further encapsulate global crowd

insights.

RAZNet and **SCAL_Net**.

4.2.2 Results and Discussions

Table 3 presents the results of different approaches on vehicle counting and traffic density estimation. Ten recently proposed approaches are compared totally. They are all effective deep learning approaches adopted from the typical computer vision task, *i.e.*, crowd counting. It can be observed that most of them perform comparably. Specifically, SFCN+ performs the best on the real part of TMS in vehicle counting task, which is a spatially fully convolutional network. RAZNet outperforms most of the others in traffic density estimation. It demonstrates the effectiveness of the zooming strategy in tiny vehicle perception once again. While CCTrans achieves the best results in terms of MAE and MSE on the game part of TMS, its performance on real part is subpar.

Besides, most approaches perform better on the game part of TMS rather than the real part, although much more training samples are provided in the real part. It is mainly because the number of vehicles varies largely in the satellite videos. Worse still, the satellite videos contain lots of noise. In contrast, the game videos are much simple. Furthermore, we can see that the performance of most approaches is boosted by conducting the augmentation strategy. It verifies the feasibility of utilizing synthesized data to make up for the lack of real satellite videos.

Fig. 7 presents the results of different approaches on the VEC and TDE task. Note that SCAL_Net performs extremely bad, so its results are not displayed. It can be observed that most of these deep learning based approaches perform well to localize the vehicles and estimate the number as well as the density. It provides an effective way to monitor the traffic from satellite.

5 CONCLUSION AND OUTLOOK

The VHR videos captured by commercial satellites enable traffic monitoring from satellite. To boost the development, we build a large-scale dataset, *i.e.*, TMS. It is composed of 12 satellite videos and 14 synthetic videos. 9296 images are obtained and fully annotated with 128,801 vehicles. It enables the application of deep learning approaches in this field and provides a new perspective to augment data from GTA-V. Three tasks are developed based on TMS, including tiny object detection, vehicle counting and density estimation. Several classic and popular approaches are tested. The results have demonstrated the challenges and superiority of TMS. The TMS dataset, annotation tools and developed tasks provide insights for traffic monitoring from satellite.

Based on the results derived from both quantitative and qualitative analysis, we have discerned captivating patterns and unveiled new challenges that demand consideration within the task of vehicle perception from satellite.

1) *How to compensate for the performance variance of different regions?* The satellite videos captured over different cities are quite different, which may cause interference to the generality of trained models and their performance vary largely. To address this problem, researchers can place particular emphasis on addressing the challenge of vehicle perception from satellite videos through domain adaptation. Furthermore, they can delve deeper into the extraction of more effective domain-invariant features from the interplay between synthetic data and real-world data.

2) *How to introduce visual foundation models to this task?* Recently, several visual foundation models are proposed, such as

SAM [57], Dino [44], GPT-4V⁹. They exceed the performance of visual perception and outperform traditional and deep learning methods significantly, which lead the datacentric research in computer vision. Researchers are encouraged to employ or tune these foundation models to the tasks of vehicle perception from satellite.

3) *How to effectively address the satellite's view displacement issue?* As mentioned in the paper, VHR videos are captured by low Earth orbit satellites. The satellite's perspective is non-stationary, allowing only a brief gaze upon the city. This renders effective vehicle tracking unfeasible. In this context, the paper did not embark on vehicle tracking tasks. Future researchers may employ techniques such as image stitching and image alignment to integrate data from multiple orbits into a coherent and real-time traffic monitoring system, enabling continuous surveillance of traffic within satellite videos.

REFERENCES

- [1] J. Zhang, X. Jia, J. Hu, and K. Tan, "Moving vehicle detection for remote sensing video surveillance with nonstationary satellite platform," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, no. 01, pp. 1–1, 2021.
- [2] Y. Shi, X. Yu, L. Liu, D. Campbell, P. Koniusz, and H. Li, "Accurate 3-dof camera geo-localization via ground-to-satellite image matching," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 2682–2697, 2023.
- [3] X. Gu, P. Angelov, C. Zhang, and P. Atkinson, "A semi-supervised deep rule-based approach for complex satellite sensor image analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [4] J. Zhang, X. Jia, and J. Hu, "Error bounded foreground and background modeling for moving object detection in satellite videos," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 58, no. 4, pp. 2659–2669, 2020.
- [5] T. Gao, Q. Niu, J. Zhang, T. Chen, S. Mei, and A. Jubair, "Global to local: A scale-aware network for remote sensing object detection," *IEEE Transactions on Geoscience and Remote Sensing*, 2023.
- [6] W. Ao, Y. Fu, X. Hou, and F. Xu, "Needles in a haystack: Tracking city-scale moving vehicles from continuously moving satellite," *IEEE Transactions on Image Processing*, vol. 29, pp. 1944–1957, 2019.
- [7] G. Heitz and D. Koller, "Learning spatial context: Using stuff to find things," in *European Conference on Computer Vision*, 2008, pp. 30–43.
- [8] C. Benedek, X. Descombes, and J. Zerubia, "Building development monitoring in multitemporal remotely sensed image pairs with stochastic birth-death dynamics," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 1, pp. 33–50, 2011.
- [9] G. Cheng, P. Zhou, and J. Han, "Learning rotation-invariant convolutional neural networks for object detection in vhr optical remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, vol. 54, no. 12, pp. 7405–7415, 2016.
- [10] Z. Liu, H. Wang, L. Weng, and Y. Yang, "Ship rotated bounding box space for ship extraction from high-resolution optical satellite images with complex backgrounds," *IEEE Geoscience and Remote Sensing Letters*, vol. 13, no. 8, pp. 1074–1078, 2016.
- [11] G.-S. Xia, X. Bai, J. Ding, Z. Zhu, S. Belongie, J. Luo, M. Datcu, M. Pelillo, and L. Zhang, "Dota: A large-scale dataset for object detection in aerial images," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3974–3983.
- [12] K. Liu and G. Mattyus, "Fast multiclass vehicle detection on aerial images," *IEEE Geoscience and Remote Sensing Letters*, vol. 12, no. 9, pp. 1938–1942, 2015.
- [13] P. F. Felzenszwalb, R. B. Girshick, and D. McAllester, "Cascade object detection with deformable part models," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2010, pp. 2241–2248.
- [14] L. Zhang and J. Ma, "Salient object detection based on progressively supervised learning for remote sensing images," *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [15] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds., 2015, pp. 91–99.

9. <https://openai.com/research/gpt-4v-system-card>

- [16] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, "SSD: single shot multibox detector," in *European Conference on Computer Vision*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., 2016, pp. 21–37.
- [17] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE Computer Society, 2016, pp. 779–788.
- [18] J. Yan, Y. Yu, X. Zhu, Z. Lei, and S. Z. Li, "Object detection by labeling superpixels," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5107–5116.
- [19] C. H. Lampert, M. B. Blaschko, and T. Hofmann, "Beyond sliding windows: Object localization by efficient subwindow search," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [20] 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.
- [21] W.-L. Zhao and C.-W. Ngo, "Flip-invariant sift for copy and object detection," *IEEE Transactions on Image Processing*, vol. 22, no. 3, pp. 980–991, 2012.
- [22] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2005, pp. 886–893.
- [23] T. Malisiewicz, A. Gupta, and A. A. Efros, "Ensemble of exemplars for object detection and beyond," in *International Conference on Computer Vision*, 2011, pp. 89–96.
- [24] P. Viola and M. Jones, "Rapid object detection using a boosted cascade of simple features," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2001, pp. 1–1.
- [25] R. Girshick, "Fast r-cnn," in *IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.
- [26] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.
- [27] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7263–7271.
- [28] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
- [29] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.
- [30] R. Guerrero-Gómez-Olmedo, B. Torre-Jiménez, R. López-Sastre, S. Maldonado-Bascón, and D. O. noro Rubio, "Extremely overlapping vehicle counting," in *Iberian Conference on Pattern Recognition and Image Analysis (IbPRIA)*, 2015.
- [31] P. Zhu, L. Wen, D. Du, X. Bian, Q. Hu, and H. Ling, "Vision meets drones: Past, present and future," *arXiv preprint arXiv:2001.06303*, 2020.
- [32] G. Kopsiaftis and K. Karantzas, "Vehicle detection and traffic density monitoring from very high resolution satellite video data," in *IEEE International Geoscience and Remote Sensing Symposium*, 2015, pp. 1881–1884.
- [33] J. Zhang, X. Jia, J. Hu, and K. Tan, "Satellite multi-vehicle tracking under inconsistent detection conditions by bilevel k-shortest paths optimization," in *Digital Image Computing: Techniques and Applications*, 2018, pp. 1–8.
- [34] "Script hook v, <http://www.dev-c.com/gtav/scripthookv/>."
- [35] E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *Journal of the ACM*, vol. 58, no. 3, pp. 1–37, 2011.
- [36] C. E. Rasmussen *et al.*, "The infinite gaussian mixture model," in *Advances in Neural Information Processing Systems*, vol. 12, 1999, pp. 554–560.
- [37] O. Barnich and M. Van Droogenbroeck, "Vibe: A universal background subtraction algorithm for video sequences," *IEEE Transactions on Image processing*, vol. 20, no. 6, pp. 1709–1724, 2010.
- [38] B. Du, Y. Sun, S. Cai, C. Wu, and Q. Du, "Object tracking in satellite videos by fusing the kernel correlation filter and the three-frame-difference algorithm," *IEEE Geoscience and Remote Sensing Letters*, vol. 15, no. 2, pp. 168–172, 2017.
- [39] X. Liu, G. Zhao, J. Yao, and C. Qi, "Background subtraction based on low-rank and structured sparse decomposition," *IEEE Transactions on Image Processing*, vol. 24, no. 8, pp. 2502–2514, 2015.
- [40] X. Zhou, C. Yang, and W. Yu, "Moving object detection by detecting contiguous outliers in the low-rank representation," *IEEE transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 3, pp. 597–610, 2012.
- [41] C. Li, L. Li, H. Jiang, K. Weng, Y. Geng, L. Li, Z. Ke, Q. Li, M. Cheng, W. Nie *et al.*, "Yolov6: A single-stage object detection framework for industrial applications," *arXiv preprint arXiv:2209.02976*, 2022.
- [42] C.-Y. Wang, A. Bochkovskiy, and H.-Y. M. Liao, "Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7464–7475.
- [43] G. Jocher, "Yolov8," *Accessed on April 05, 2023*.
- [44] H. Zhang, F. Li, S. Liu, L. Zhang, H. Su, J. Zhu, L. Ni, and H.-Y. Shum, "DINO: DETR with improved denoising anchor boxes for end-to-end object detection," in *The Eleventh International Conference on Learning Representations*, 2023. [Online]. Available: <https://openreview.net/forum?id=3mRwyG5one>
- [45] C. Feng, Y. Zhong, Y. Gao, M. R. Scott, and W. Huang, "Tood: Task-aligned one-stage object detection," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. Los Alamitos, CA, USA: IEEE Computer Society, oct 2021, pp. 3490–3499. [Online]. Available: <https://doi.ieeecomputersociety.org/10.1109/ICCV48922.2021.00349>
- [46] Y. Wang, X. Hou, and L.-P. Chau, "Dense point prediction: A simple baseline for crowd counting and localization," in *IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, 2021, pp. 1–6.
- [47] C. Liu, X. Weng, and Y. Mu, "Recurrent attentive zooming for joint crowd counting and precise localization," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1217–1226.
- [48] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma, "Single-image crowd counting via multi-column convolutional neural network," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 589–597.
- [49] J. Gao, W. Lin, B. Zhao, D. Wang, C. Gao, and J. Wen, "C³ framework: An open-source pytorch code for crowd counting," *arXiv preprint arXiv:1907.02724*, 2019.
- [50] Y. Li, X. Zhang, and D. Chen, "Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1091–1100.
- [51] Q. Wang, J. Gao, W. Lin, and Y. Yuan, "Learning from synthetic data for crowd counting in the wild," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8198–8207.
- [52] J. Gao, Q. Wang, and Y. Yuan, "Scar: Spatial/channel-wise attention regression networks for crowd counting," *Neurocomputing*, vol. 363, pp. 1–8, 2019.
- [53] J. Wan, Z. Liu, and A. B. Chan, "A generalized loss function for crowd counting and localization," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 1974–1983.
- [54] Q. Song, C. Wang, Z. Jiang, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, and Y. Wu, "Rethinking counting and localization in crowds: A purely point-based framework," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 3365–3374.
- [55] Y. Tian, X. Chu, and H. Wang, "Cctrans: Simplifying and improving crowd counting with transformer," *arXiv preprint arXiv:2109.14483*, 2021.
- [56] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [57] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo *et al.*, "Segment anything," *arXiv preprint arXiv:2304.02643*, 2023.