

# TJU-DHD: A Diverse High-Resolution Dataset for Object Detection

Yanwei Pang, Jiale Cao, Yazhao Li, Jin Xie, Hanqing Sun, Jinfeng Gong

**Abstract**—Vehicles, pedestrians, and riders are the most important and interesting objects for the perception modules of self-driving vehicles and video surveillance. However, the state-of-the-art performance of detecting such important objects (esp. small objects) is far from satisfying the demand of practical systems. Large-scale, rich-diversity, and high-resolution datasets play an important role in developing better object detection methods to satisfy the demand. Existing public large-scale datasets such as MS COCO collected from websites do not focus on the specific scenarios. Moreover, the popular datasets (e.g., KITTI and Citypersons) collected from the specific scenarios are limited in the number of images and instances, the resolution, and the diversity. To attempt to solve the problem, we build a diverse high-resolution dataset (called TJU-DHD). The dataset contains 115,354 high-resolution images (52% images have a resolution of  $1624 \times 1200$  pixels and 48% images have a resolution of at least  $2,560 \times 1,440$  pixels) and 709,330 labeled objects in total with a large variance in scale and appearance. Meanwhile, the dataset has a rich diversity in season variance, illumination variance, and weather variance. In addition, a new diverse pedestrian dataset is further built. With the four different detectors (i.e., the one-stage RetinaNet, anchor-free FCOS, two-stage FPN, and Cascade R-CNN), experiments about object detection and pedestrian detection are conducted. We hope that the newly built dataset can help promote the research on object detection and pedestrian detection in these two scenes. The dataset is available at <https://github.com/tjuiit/TJU-DHD>.

**Index Terms**—Dataset, object detection, pedestrian detection, high resolution, large scale.

## I. INTRODUCTION

Object detection aims to locate and classify objects in an image, which is a fundamental but challenging task in computer vision community. In recent few years, based on deep Convolutional Neural Networks (CNN) [27], [55], [29], [31], object detection has achieved great progress and started to be successfully applied to real life. Behind the technique of deep CNN, the large-scale image datasets, such as ImageNet [53], PASCAL VOC [18], and MS COCO [39], are another key to push the progress of object detection. These large-scale datasets (e.g., MS COCO [39]) collected from the website do not focus on any specific scene. As a result, the detector trained on these about generic object datasets cannot achieve

The work is supported by the National Key R&D Program of China (Grant No. 2018AAA0102800 and 2018AAA0102802) and National Natural Science Foundation of China (Grant No. 61632018 and 61906131). (The corresponding author: Jiale Cao)

Y. Pang, J. Cao, Y. Li, J. Xie, and H. Sun are with the School of Electrical and Information Engineering and Tianjin Key Laboratory of Brain-inspired Intelligence Technology, Tianjin University, Tianjin 300072, China (E-mail: {pyw, connor, lyztju, jinxie, hqSun}@tju.edu.cn).

J. Gong is with the China Automotive Technology and Research Center Co., Ltd., Tianjin 3003000, China (E-mail: gongjinfeng@catarc.ac.cn).

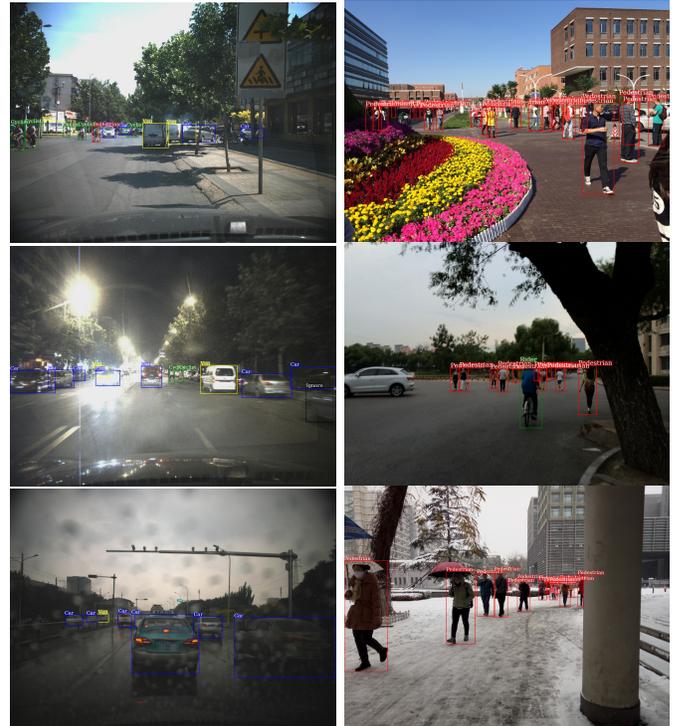


Fig. 1. Some examples of the newly built TJU-DHD. The dataset is collected under different scenes, different weathers, different seasons, and different illuminations. The high-resolution images in the first column are from TJU-DHD-traffic and those in the second column are from TJU-DHD-campus.

a very good performance in the specific application scene. It is necessary to develop a large-scale object dataset in specific scenes for specific application requirements.

In our daily life, traffic scene and campus scene are two common scenes. It is important to detect objects (e.g., pedestrian and car) in these two scenes for real applications. Some datasets [15], [20], [75] in these specific applications have been built in the past decade. For example, in the urban traffic scene, Geiger *et al.* [20] built the KITTI benchmark for object detection, and Dollár *et al.* [15] proposed the Caltech dataset for pedestrian detection. The KITTI dataset [20] contains 7,481 training images and 7,518 test images with a resolution of  $1,240 \times 376$  pixels, while the Caltech dataset [15] has a low resolution of  $640 \times 480$  pixels and limited pedestrians. We argue that these two datasets are still not enough to push the progress of object detection in the specific scenes. Recently, some large-scale datasets (e.g., ApolloScape dataset [32] and BDD100K [73]) have been proposed. However, the ApolloScape dataset [32] does not cover the season variance

and day-night variance, and the BDD100K dataset [73] has a relatively low resolution of 720p. Meanwhile, they merely focus on the traffic scene.

In this paper, we develop a new Diverse High-resolution Dataset (called TJU-DHD, collected by Tianjin university) in the traffic scene and campus scene. Fig. 1 shows some examples. It contains 115,354 images and 709,330 labeled instances. The resolution of the images is of at least  $1,624 \times 1,200$  pixels and the height of the objects ranges from 11 pixels to 4,152 pixels. Meanwhile, the dataset collected over one year has a large variance in object appearance, object scale, illumination, season, and weather. As a result, the dataset has a very rich diversity. For object detection, TJU-DHD is split into two subsets: TJU-DHD-traffic subset in the traffic scene and TJU-DHD-campus subset in the campus scene. To further focus on one of the most important cases of object detection (*i.e.*, pedestrian detection), we choose the images including the pedestrians in these two scenes to construct a large-scale pedestrian dataset. Based on the built dataset (*i.e.*, TJU-DHD), we implement the one-stage method RetinaNet [41], the anchor-free method FCOS [59], the two-stage method FPN [40], and the cascade method Cascade R-CNN [5] to give the baseline performance. To summarize, the contributions of this paper can be summarized as follows.

(1) A new diverse high-resolution dataset for object detection in two important scenes is built, which has a rich variance in appearance, scale, illumination, season, and weather.

(2) A new large-scale pedestrian dataset is further proposed based on the built diverse high-resolution object dataset, which can provide both the same-scene and cross-scene evaluations.

(3) Experiments based on four different detectors (*i.e.*, one-stage RetinaNet [41], anchor-free FCOS [59], two-stage FPN [40], and Cascade R-CNN [5]) are conducted to provide the baseline performance. We hope that this newly built high-resolution dataset can push the progress of object detection and pedestrian detection in the traffic scene and campus scene.

## II. RELATED WORK

In this section, we firstly give a brief review of object detection datasets, including pedestrian detection datasets and face detection datasets, and secondly give a review of object detection methods, including pedestrian detection methods.

### A. The datasets of object detection

To promote the progress of object detection and provide a fair performance comparison, many datasets for object detection have been proposed in the past decade. In this subsection, we review some object detection datasets [18], [39], [20], some pedestrian detection datasets [15], [75], [54], and some face detection datasets [70], [85], [34].

*Generic object detection datasets* Generally, these object datasets are usually collected from the website and do not focus on the specific application scenes. PASCAL VOC [18] and MS COCO [39] are two of the most famous generic object datasets. The PASCAL VOC challenge has been held since 2006. Among these PASCAL VOC datasets, VOC2007 and VOC2012 are widely used, which have 20 object classes

and over 11,000 images. Compared with PASCAL VOC, MS COCO [39] is significantly larger in category and images, which contains 80 object classes and about 328k images. Recently, some newly built datasets (*e.g.*, LVIS [23] and OpenImages [36]) contain a very large number of object classes and images. The LVIS dataset has 1,000 object classes and 164k images, and the OpenImages dataset has 600 object classes and 9.2M images. Besides these generic object datasets, Some datasets, including KITTI [20], Cityscapes [12], Mapillary vistas [48], ApolloScape [32], and BDD100K [73], focus on the traffic scene. Specifically, the KITTI dataset has 3 classes and 14,999 images, the Cityscapes benchmark has 20 classes and 5,000 fine-annotation images, the Mapillary Vistas dataset has 66 classes and 25,000 images, and the BDD100K dataset has 10 classes and 100k images.

*Pedestrian detection datasets* Pedestrian detection is a key task in both self-driving and video surveillance [20], [68], [67], [69]. The INRIA [13], ETH [16], and Daimler [17] datasets are the early datasets of pedestrian detection. After those datasets, Dollár *et al.* [15] built the larger Caltech pedestrian dataset and gave a unified evaluation of pedestrian detection. The standard Caltech pedestrian dataset consists of 4,250 images for training and 4,024 for test. However, there are only 0.32 pedestrians per image and the image resolution is low (*i.e.*,  $640 \times 480$  pixels). Recently, a new pedestrian dataset Citypersons [75] extracted from Cityscapes is proposed. There are 6.5 pedestrians per image and the image resolution is of  $2,048 \times 1,024$  pixels. Similar to Citypersons, another new large-scale pedestrian dataset EuroCity [1] is also collected in multiple European cities. To promote the progress of pedestrian detection in crowd scenes, a new dataset CrowdHuman [54] is collected from the website, which contains about 25,000 images.

*Face detection datasets* Face detection is another classic task in computer vision, which is important for face recognition and face verification. The common face detection datasets contain AFW [85], FDDB [34], PASCAL FACE [71], and WiderFace [70]. The AFW dataset [85] has 205 images with 473 faces. The FDDB dataset [34] has 2,845 images with 5,171 faces. The PASCAL FACE dataset [34] has 851 images with 1,341 faces. However, these datasets have the limited images and faces. To better cover face variations in pose, scale, and occlusion, the WiderFace dataset [70] is proposed, which has 32,203 images with 39,3703 faces.

### B. The methods of object detection

In the past decade, the dominant methods for object detection have changed from handcrafted features based ones to deep Convolutional Neural Networks (CNN) based ones. Before the CNN based methods, the researchers proposed many widely used handcrafted features (*e.g.*, Haar [61], LBP [50], [62], ICF [14], and HOG [13]) for object detection. With the success of deep CNN on image classification [27], [55], [29], [31], the researchers started to apply deep CNN to promote the progress of object detection. The CNN based methods can be mainly divided into two different classes: two-stage methods [21], [22], [25], [30], [10] and one-stage methods [43], [24], [41], [26], [49].

Two-stage methods firstly extract some candidate object proposals and secondly classify these proposals into specific object classes. R-CNN [21] is the first two-stage method, which firstly uses the handcrafted features based method selective search [60] for proposal extraction, secondly calculates the CNN features for each candidate proposal, and finally classifies these proposals into specific object classes by SVM. To reduce computational costs, Fast R-CNN [22] and SPPnet [28] share the CNN feature calculations of all the candidate proposals by firstly generating the CNN features of whole image. Based on Fast R-CNN, Faster R-CNN [25] further joins proposal extraction and proposal classification in an end-to-end network. Faster R-CNN [25] not only improves the quality of proposal extraction but also reduces computational costs of proposal extraction. To solve the scale-variance problem in object detection, both feature pyramid methods [40], [4], [10], [38] and image pyramid methods [56], [57], [64], [47] are proposed. Among the feature pyramid methods, MS-CNN [4] and FPN [40] are two representative methods, which use the shallow and high-resolution layer for small-scale object detection and use the deep and low-resolution layer for large-scale object detection. Generally, the image pyramid methods rescale the input image into a sequence of images at different scales and detect objects on each re-scaled image. To tackle the wide scale spectrum at the training stage, SNIP [56] selectively back-propagates the gradients of objects at different sizes as a function of the image scale. Recently, Mask R-CNN [30] is proposed to extend object detection to instance segmentation by an additional branch for mask prediction.

One-stage methods directly predict object classes and box regression offsets of dense boxes. YOLO [24] and SSD [43] are two representative one-stage methods. YOLO [24] splits the input image into  $N \times N$  grids and predicts object probability and class in each grid. Similar to FPN [40], SSD [43] uses the shallow layer to detect small-scale objects and uses the deep layer to detect large-scale objects. Compared with two-stage methods, one-stage methods suffer more from class imbalance problem. To solve this problem, focal loss [41] is proposed to down-weight the easy examples and up-weight the hard examples. To enhance the semantic of the features, many improvements [80], [79], [6] have been proposed. To improve location precision, some methods [77], [35], [9] use the cascade structure to perform the regression more than once. The above one-stage methods are anchor-based methods, which needs to set some hyper-parameters (*e.g.*, anchor scales and aspect ratios). To alleviate the drawbacks of the empirical hyper-parameters introduced by the anchor-based methods, some anchor-free methods [37], [81], [59], [72], [19] are recently proposed. CornerNet [37] uses the top-left and bottom-right corners to locate and classify objects. CenterNet [81] predicts the center points, the heights, and the widths of objects. FCOS [59] uses a feature pyramid structure for anchor-free object detection.

As a special and important case of object detection, pedestrian detection [14], [84], [76], [65], [44], [82] has also attracted much attention of researchers. Compared with generic object detection, pedestrian detection faces more severe challenges in scale variance and occlusion. Before the deep

TABLE I  
THE NUMBER OF IMAGES AND INSTANCES IN THE TRAINING SET, THE VALIDATION SET, AND THE TEST SET OF THE NEW BUILT DATASET (CALLED TJU-DHD). TJU-DHD CONTAINS TWO SUBSETS, TJU-DHD-TRAFFIC AND TJU-DHD-CAMPUS.

Name	TJU-DHD-traffic		TJU-DHD-campus	
	#images	#instances	#images	#instances
training set	45,266	239,980	39,727	267,445
validation set	5,000	30,679	5,204	41,620
test set	10,000	60,963	10,157	68,643
total	60,266	331,622	55,088	377,708



Fig. 2. Examples of annotated objects in five classes (*i.e.*, car, van, truck, rider, and pedestrian) in the built TJU-DHD. In our dataset, sedan and SUV are treated as the same class (called car), van, minibus, and bus are all treated as the same class (called van), and cyclist, motorcyclist, and tricyclist are all treated as the same class (called rider).

CNN based methods, the handcrafted channel features based methods [14], [74], [7] are dominant, which first convert the color images to ten channel images (*i.e.*, three LUV channels, one gradient magnitude channel, and six oriented gradient channels) and second extract local and non-local features to learn a pedestrian detector. Recently, deep CNN based methods greatly push the progress of pedestrian detection. Some methods [46], [3], [2], [42] use semantic information to improve pedestrian detection. Some methods [58], [8], [52], [66] aim to improve small-scale pedestrian detection, while some methods [78], [83], [51] exploit the part or visible information for occluded pedestrian detection. To improve pedestrian detection in crowded scenes, some methods [63], [45], [33], [11] exploit how to combine the highly overlapping bounding boxes.

### III. OUR DATASET DETAILS

In this section, we firstly introduce the newly built diverse high-resolution dataset (called TJU-DHD) under traffic scene and campus scene, secondly describe the new large-scale pedestrian dataset under these two scenes, and finally compare the built dataset with some related datasets.

#### A. Images, object annotations, and dataset splits

Traffic scene and campus scene are two common scenes in our daily life. Due to the complexity of the scenes, detecting objects in these two scenes are relatively difficult. To promote the progress of object detection in these two scenes, we build a new diverse high-resolution dataset, including two different subsets corresponding to these two scenes (called TJU-DHD-traffic and TJU-DHD-campus).

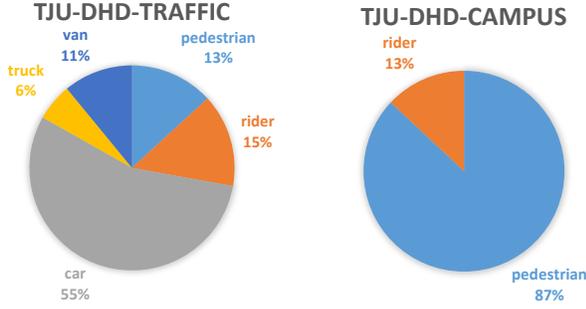


Fig. 3. The proportions of different object classes.

TJU-DHD-traffic subset is collected by a driving car in the traffic scene, which has 60,266 images and 331,622 labeled instances (see Table I). The images have a fixed high resolution (*i.e.*,  $1,624 \times 1,200$  pixels). We select five common and important object classes, namely car, van, truck, pedestrian, and rider, for bounding box annotations. Fig. 2 gives some examples and illustrations of five different object classes. To annotate each object instance, a bounding box  $(x_1, y_1, x_2, y_2)$  represented by the left-top point  $(x_1, y_1)$  and the right-bottom point  $(x_2, y_2)$  is used. Meanwhile, the occlusion level and the truncation level of each instance are also given. The TJU-DHD-traffic is split into three sets: the training set, the validation set, and the test set. Table I shows the number of images and instances in each set. Specifically, the training set contains 45,266 images and 239,980 instances, the validation set contains 5,000 images and 30,679 instances, and the test set contains 10,000 images and 60,963 instances.

TJU-DHD-campus subset is taken from multiple different mobile phones mainly on the university campus, which has 55,088 images and 377,708 labeled instances. Because the images are taken by different mobile phones, the image resolutions are high but various, which are of at least  $2,560 \times 1,440$  pixels. We only select two most important and common objects (*i.e.*, pedestrian and rider) in campus scene for bounding box annotations. To better exploit pedestrian detection for the researchers, especially occluded pedestrian detection, we provide two bounding box annotations for each pedestrian. The two bounding boxes respectively represent the full body and the visible part of the pedestrian. The TJU-DHD-campus is also split into three different sets: the training set, the validation set, and the test set. Table I shows the number of images and instances in each set. The training set contains 39,727 images and 267,445 instances, the validation set contains 5,204 images and 41,620 instances, and the test set contains 10,157 images and 68,643 instances.

### B. Dataset statistics

Based on the newly built dataset (called TJU-DHD), we give the detailed statistics about category variance, scale variance, occlusion variance, and dataset diversity to better know and understand this newly built dataset.

**Category statistics** Fig. 3 shows the proportions of different object classes in both the TJU-DHD-traffic and TJU-DHD-campus. In the TJU-DHD-traffic, the proportions of pedestrian,

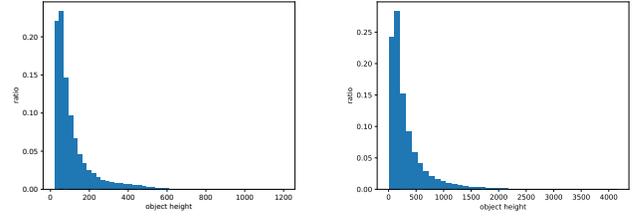


Fig. 4. Scale variance in the TJU-DHD-traffic (left) and TJU-DHD-campus (right).

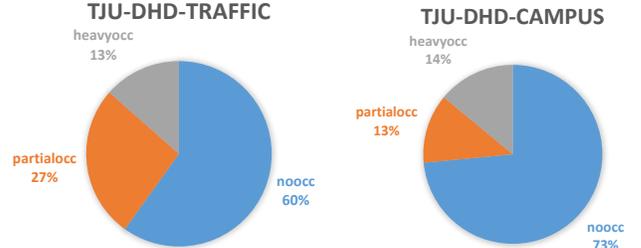


Fig. 5. The proportions of different occlusion levels.

rider, car, truck, and van are 13%, 15%, 55%, 6%, and 11%, respectively. In the TJU-DHD-campus, the proportions of pedestrian and rider are 87% and 13%, respectively.

**Scale statistics** Fig. 4 shows the scale variance of objects in both the TJU-DHD-traffic and TJU-DHD-campus. we use the object height to represent its scale because the height and the scale are closely related under these two scenes. Fig. 4(a) shows the object height variance in the TJU-DHD-traffic. The heights vary from 20 pixels to 1,200 pixels and most objects have a height of fewer than 200 pixels. Fig. 4(b) shows the height variance in the TJU-DHD-campus, where the height varies from 11 pixels to 4,152 pixels. Namely, our dataset has a very large variance in object scale.

**Occlusion statistics** Similar to the Caltech and Citypersons datasets [15], [75], object occlusion level is divided into three different levels (*i.e.*, no occlusion, partial occlusion, and heavy occlusion). Partial occlusion means that the occlusion ratio  $o$  is less than 0.35 and greater than 0.0 (*i.e.*,  $o \leq 0.35$ ), and heavy occlusion means that the occlusion ratio  $o$  is greater than 0.35 (*i.e.*,  $o > 0.35$ ). Fig. 5 plots the proportions of different object occlusion levels in both the TJU-DHD-traffic and TJU-DHD-campus. In the TJU-DHD-traffic, the proportion of no occlusion is 60%, the proportion of partial occlusion is 27%, and the proportion of heavy occlusion is 13%. In the TJU-DHD-campus, the proportion of no occlusion is 73%, the proportion of partial occlusion is 13%, and the proportion of heavy occlusion is 14%.

**Diversity** The datasets are collected over more than one year. Except for the diversity in object appearance, object scale, and object density, it also covers the diversity in illumination variance, scene variance, weather variance, and season variance. Fig. 6(a) gives some examples of illuminate variance. Because the built dataset is captured from day to night, the illumination variance can be caused by daytime, nighttime, frontlight, and backlight. Fig. 6(b) gives some examples of traffic scene



Fig. 6. The rich diversity of the built TJU-DHD, which contains the variances of illumination, scene, weather, and season.

TABLE II  
THE NUMBER OF IMAGES AND INSTANCES IN THE TRAINING SET, THE VALIDATION SET, AND THE TEST SET FOR PEDESTRIAN DETECTION UNDER TWO SCENES.

Name	TJU-Ped-traffic		TJU-Ped-campus	
	#images	#instances	#images	#instances
training set	13,858	27,650	39,727	234,455
validation set	2,136	5,244	5,204	36,161
test set	4,344	10,724	10,157	59,007
total	20,338	43,618	55,088	329,623

variance, which contains urban road, highway road, and rural road. Fig. 6(c) gives some examples of weather variance, which contains sunny days, cloudy days, and rainy days. Fig. 6(d) gives some examples of season variance from spring to winter. The season variance can cause a large variance in the appearances of both objects and background.

### C. The newly built pedestrian dataset

Pedestrian detection is a typical and important case in object detection, which has drawn much attention of the researchers in the past decade. To better focus on the specific pedestrian detection, a new large-scale pedestrian dataset is built based on the TJU-DHD-traffic and TJU-DHD-campus. To simplify the expression, the newly built large-scale pedestrian dataset is called TJU-DHD-pedestrian.

We choose the images which contain the pedestrians to construct TJU-DHD-pedestrian. In TJU-DHD-pedestrian, the

annotations of car, van, and truck are removed. Because the rider has a relation to the pedestrian, the annotations of rider are set as the ignored regions. As a result, TJU-DHD-pedestrian has 75,426 images, 373,241 labeled pedestrians, and 112,842 ignored bounding boxes (see Table II). Instead of combining the images in two scenes together, we keep the original split of the training, validation, and test sets. Thus, TJU-DHD-pedestrian has two training, validation, and test sets. Specifically, there are 20,338 images and 43,618 labeled pedestrians in the traffic scene (called TJU-Ped-traffic), and there are 55,088 images and 329,623 labeled pedestrians in the campus scene (called TJU-Ped-campus). These two sets have a large gap in image resolution, detection scene, and instance scale. Based on these two sets, we can give a better performance analysis of the pedestrian detector with both the same-scene evaluation and the cross-scene evaluation.

### D. Comparison of some related datasets

In this subsection, some related object datasets (*e.g.*, Caltech [15], Citypersons [75], CrowdedHuman [54], EuroCity [1], and KITTI [20]) are compared with the newly built dataset (TJU-DHD) in Table III. For a fair comparison, the statistics are only performed on the training set.

*Caltech*<sup>1</sup> is a very popular pedestrian dataset in the past decade [15]. The standard training set has 4,250 images. To improve detection performance, some researchers [74], [7]

<sup>1</sup>[http://www.vision.caltech.edu/Image\\_Datasets/CaltechPedestrians](http://www.vision.caltech.edu/Image_Datasets/CaltechPedestrians)

TABLE III

COMPARISONS WITH SOME RELATED DATASETS. FOR A FAIR COMPARISON, THE STATISTICS ARE ONLY BASED ON THE TRAINING SET. ‘#AVG’, ‘VBOX’, AND ‘#IGNORE’ REPRESENT AVERAGED OBJECT NUMBER PER IMAGE, VISIBLE BOX ANNOTATION, AND THE NUMBER OF IGNORED BOUNDING BOXES.

Name	scene	day/night	seasons	#images	resolution	#objects	#avg	vbox	categories	#ignore
Caltech [15]	traffic	day	1	42,782	640 × 480	13,674	0.32	✓	ped.	51,092
KITTI [20]	traffic	day	1	7,481	1,240 × 376	40,570	5.42	✗	ped., car, cyclist	11,295
Citypersons [75]	traffic	day	1	2,975	2,048 × 1,024	19,238	6.47	✓	ped.	6,768
CrowdedHuman [54]	crowd	day	1	15,000	> 400 × 300	339,565	22.64	✓	people	99,227
EuroCity [1]	traffic	day,night	4	28,114	1,920 × 1,024	142,736	5.08	✗	ped., rider	83,218
TJU-DHD-traffic	traffic	day,night	1	45,266	1,624 × 1,200	239,980	5.30	✗	ped., rider, car, truck, van	11,711
TJU-DHD-campus	campus	day,night	4	39,727	> 2,560 × 1,440	267,445	6.73	✓	ped., rider	24,321
TJU-DHD-pedestrian	traffic, campus	day,night	4	53,585	≥ 1,624 × 1,200	262,105	4.89	✓	ped.	73,846

started to enlarge the training set by capturing every third frame. As a result, it has a large number of images (*i.e.*, 42,782). However, two adjacent images are highly correlated and there are only 13,674 pedestrians. Meanwhile, the image resolution of 640×480 pixels is relatively low.

*KITTI*<sup>2</sup> is a challenging benchmark [20] for the application of self-driving, which contains many different computer vision tasks (*i.e.*, stereo, optical flow, visual odometry, object detection, and tracking). For object detection task, the training set has 7,481 images with a resolution of 1,240×376 pixels. Meanwhile, there are 19,238 labeled objects. Thus, the numbers of images and the instances are both limited.

*Citypersons*<sup>3</sup> is a recently built pedestrian dataset [75] collected from the Cityscapes benchmark [12]. Compared with the Caltech and KITTI datasets [15], [20], Citypersons has more pedestrians per image and larger image resolution (*i.e.*, 2,048×1,024 pixels). However, the numbers of images and instances are still limited.

*CrowdedHuman*<sup>4</sup> is collected from the website and aims to improve detection performance in the human crowd scene [54]. It does not focus on the specific application scene (*e.g.*, the traffic scene or the campus scene). Meanwhile, the image resolution in CrowdedHuman is relatively low.

*EuroCity persons*<sup>5</sup> is collected from multiple cities in Europe [1]. Compared with the previous datasets [15], [75], EuroCity persons has a larger number of images (*i.e.*, 40,219) under different illuminations from day to night. Compared with Eurocity, our proposed TJU-DHD has the following differences: (1) Our dataset is more complete, which focuses on not only pedestrian/rider detection but also vehicle detection. (2) Our pedestrian dataset is almost two times larger than EuroCity in both images and objects (see Table III). Moreover, it provides visible box annotation and two different scenes (traffic and campus). (3) Our pedestrian dataset is far from saturated. In Table X, the miss rate on the R set of EuroCity is 6.81%, while miss rates on the R set of TJU-DHD-traffic and TJU-DHD-campus are 27.92% and 22.30%. Thus, our pedestrian dataset has more space for future research.

Overall, the advantages of our built dataset can be summarized as follows: (1) *A large amount of data.* TJU-DHD-

traffic, TJU-DHD-campus, TJU-DHD-pedestrian almost have the largest number of both images and instances. For example, TJU-DHD-pedestrian has almost two times larger than EuroCity in both images and pedestrians. (2) *High resolution.* Our built dataset has higher resolutions than other dataset, which can provide more space for the research of small-scale object detection. TJU-DHD-traffic has a fixed high resolution of 1,624×1,200 pixels, while TJU-DHD-campus has a high resolution of at least 2,560×1,400 pixels. (3) *Large time span.* On the one hand, TJU-DHD-traffic and TJU-DHD-campus are collected from day to night. On the other hand, TJU-DHD-campus is collected over one year. Thus, it has a rich diversity. (4) *Cross-scene evaluation.* TJU-DHD-pedestrian contains the pedestrians under two different scenes, which can not only give the same-scene evaluation but also give the cross-scene evaluation.

## IV. EXPERIMENTS

Recently, deep convolutional neural networks [21], [22], [25] have achieved great success in object detection. In this section, we select four representative methods (*i.e.*, RetinaNet [41], FCOS [59], FPN [40], and Cascade R-CNN [5]) to conduct some fundamental experiments on the newly built dataset and give the baseline for future research.

### A. Evaluation metric

Following the standard evaluation on the MS COCO benchmark [39] and Caltech pedestrian dataset [15], mean average precision (mAP) and log-average miss rate (MR) are respectively used to evaluate object detection and pedestrian detection, which are introduced as follows.

*Mean average precision* The average precision (AP) on the MS COCO benchmark is averaged under ten different intersection over union (IoU) thresholds of 0.50:0.05:0.95. Compared with AP on the PASCAL VOC [18], mAP on the MS COCO considers not only the accuracy of object classification but also the precision of object location.

*Log-average miss rate* It is computed by averaging log miss rates at nine different FPPI rates varying from  $10^{-2}$  to  $10^0$ . FPPI means false positive per image. The IoU threshold used to assign a box as positive or negative is set to be 0.5.

To evaluate the detection performance under different scales, the objects are divided into different sets according to the

<sup>2</sup><http://www.cvlibs.net/datasets/kitti>

<sup>3</sup><https://bitbucket.org/shanshanzhang/citypersons>

<sup>4</sup><https://www.crowdhuman.org>

<sup>5</sup><https://eurocity-dataset.tudelft.nl>

TABLE IV  
AVERAGE PRECISIONS (AP) OF SOME STATE-OF-THE-ART DETECTORS ON THE TJU-DHD-TRAFFIC.

method	backbone	input size	AP	AP@0.5	AP@0.75	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>
(a) RetinaNet [41]	ResNet50	1,333×800	53.5	80.9	60.0	24.0	50.5	68.0
(b) FCOS [59]	ResNet50	1,333×800	53.8	80.0	60.1	24.6	50.6	68.8
(c) FPN [40]	ResNet50	1,333×800	55.4	83.4	63.0	30.4	52.2	68.2
(d) Cascade R-CNN [5]	ResNet50	1333×800	57.9	82.7	66.6	32.6	54.4	71.4

object scales. In the TJU-DHD-traffic, following the same protocol as used in the MS COCO benchmark, we divide the objects into three different sets (*i.e.*, small objects: area  $< 32^2$ , middle objects:  $32^2 < \text{area} < 96^2$ , large objects: area  $> 96^2$ ). In the TJU-DHD-campus, due to the much larger scale variations compared with the COCO benchmark, we divide the objects into four different sets according to the heights of the objects, namely tiny objects (height  $< 80$ ), small objects ( $80 < \text{height} < 160$ ), medium objects ( $160 < \text{height} < 320$ ), and large objects (height  $> 320$ ). In TJU-DHD-campus dataset, there are 14.0% tiny objects, 26.8% small objects, 26.6% medium objects, and 32.6% large objects.

### B. Baseline detectors

Because the feature pyramid structure is very useful for multi-scale object detection and pedestrian detection, we select four representative feature pyramid detectors to conduct the experiments on the newly built dataset (called TJU-DHD). They are the one-stage detector RetinaNet [41], the anchor-free detector FCOS [59], the two-stage detector FPN [40], and the cascade detector Cascade R-CNN [5].

*RetinaNet* To solve the class imbalance problem in single-stage methods (*e.g.*, YOLO [24] and SSD [43]) and improve detection accuracy, RetinaNet [41] adopts focal loss to decrease the weights of easy samples and increase the weights of hard samples. As a result, it achieves a comparable performance with two-stage methods.

*FCOS* RetinaNet [41] is an anchor-based method, which is required to design the scales and aspect ratios of anchors. Compared to RetinaNet, FCOS [59] is an anchor-free method, which predicts the offsets of left, right, top, and down for each point inside an object bounding box. Thus, FCOS does not require the handcrafted design of anchors.

*FPN* Compared to Faster R-CNN [25], FPN [40] uses the in-network layers of different resolutions to detect the objects at different scales and adopts the top-down structure to enhance the feature semantic of each output layer. Thus, FPN can largely improve the performance of object detection, especially small-scale object detection.

*Cascade R-CNN* To improve the localization accuracy, Cascade R-CNN [5] stacks multiple ROI detectors. The different ROI detectors are trained stage by stage with increasing IoU thresholds. As a result, Cascade R-CNN can progressively improve object localization accuracy.

### C. Implementation details

These four detectors are implemented based on the open source object detection toolbox mmdetection<sup>6</sup>. The widely

<sup>6</sup><https://github.com/open-mmlab/mmdetection>

TABLE V  
AVERAGE PRECISIONS (AP) PER CATEGORY ON THE TJU-DHD-TRAFFIC.

method	AP	AP <sub>ped</sub>	AP <sub>rid</sub>	AP <sub>car</sub>	AP <sub>tru</sub>	AP <sub>van</sub>
(a) RetinaNet [41]	53.5	58.4	35.9	48.0	69.7	55.7
(b) FCOS [59]	53.8	58.7	45.0	48.4	71.2	56.0
(c) FPN [40]	55.4	60.0	39.7	50.6	70.9	55.7
(d) Cascade R-CNN [5]	57.9	61.7	42.5	53.7	73.0	58.7

used deep residual model ResNet50 [29] is used as the backbone. In the object detection experiments, the typical input with the resolution of 1,333×800 pixels on the COCO benchmark is used for training and testing. In the pedestrian detection experiments, the typical input with the resolution of 2,048×1,024 pixels on the Citypersons is used for training and testing. We use 2 NVIDIA GPUs for training. A mini-batch has 4 images per GPUs. The initial learning rate is 0.01 for FPN, 0.005 for RetinaNet, 0.005 for FCOS, and 0.01 for Cascade R-CNN. After that, the learning rate decreases at epoch 8 and epoch 11 by a factor of 10. During the inference stage, the top 100 detection bounding boxes per image are saved for performance evaluation.

### D. Experimental results on the TJU-DHD-traffic

In this subsection, the experiments on TJU-DHD-traffic are conducted. Table IV shows the average precisions (AP) of these four detectors (*i.e.*, RetinaNet [41], FCOS [59], FPN [40], and Cascade R-CNN [5]). It can be concluded as follows: (1) The anchor-free FCOS is a little superior (+0.3%) to the anchor-based RetinaNet and inferior (-1.6%) to the two-stage FPN. (2) The two-stage FPN is 1.9% better than one-stage RetinaNet. On small-scale object detection, FPN is 6.4% better than RetinaNet. (3) Cascade structure is very useful for accurate object detection. For example, Cascade R-CNN is even 0.7% worse than FPN on AP@0.5, while it is 3.6% better than FPN on AP@0.75.

To see the performance on each object category, Table V further shows the AP per object category. Cascade R-CNN steadily outperforms the other three methods on each category, which further shows the effectiveness of cascade structure. The four detectors all have a relatively lower performance on the categories of car and rider. The reason may be explained as follows: (1) riders have a large inter-class variance because it contains cyclist, motorcyclist, and tricyclist. (2) Most cars in our traffic scene are very crowded and occluded, which are relatively difficult for detection.

TABLE VI  
AVERAGE PRECISIONS (AP) OF SOME STATE-OF-THE-ART DETECTORS ON THE TJU-DHD-CAMPUS.

method	backbone	input size	AP	AP@0.5	AP@0.75	AP <sub>t</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>
(a) RetinaNet [41]	ResNet50	1,333×800	48.4	73.9	52.4	4.7	27.3	56.2	73.8
(b) FCOS [59]	ResNet50	1,333×800	49.3	73.8	53.8	5.6	29.6	55.9	74.3
(c) FPN [40]	ResNet50	1,333×800	52.4	77.5	58.4	8.5	37.4	58.6	74.9
(d) Cascade R-CNN [5]	ResNet50	1,333×800	55.1	77.6	60.9	10.8	40.1	61.2	78.0

TABLE VII  
AVERAGE PRECISIONS (AP) PER CATEGORY ON THE TJU-DHD-CAMPUS.

method	AP	AP <sub>ped</sub>	AP <sub>rider</sub>
(a) RetinaNet [41]	48.4	50.5	46.4
(b) FCOS [59]	49.3	51.1	47.5
(c) FPN [40]	52.4	54.4	50.4
(d) Cascade R-CNN [5]	55.1	57.2	53.1

TABLE VIII  
IMPACT OF DIFFERENT INPUT SIZES (RESOLUTIONS) ON DETECTION PERFORMANCE ON THE TJU-DHD-CAMPUS. ‘ $N \times \downarrow$  & ‘ $\uparrow$ ’ INDICATES FIRST  $N \times$  DOWNSAMPLING THE IMAGES AND SECOND  $N \times$  UPSAMPLING THE IMAGES. THE RUN-TIME (MS) IS TESTED ON A SINGLE NVIDIA P6000 GPU.

input size	AP	AP <sub>t</sub>	AP <sub>s</sub>	AP <sub>m</sub>	AP <sub>l</sub>	time
1333×800	52.4	8.5	37.4	58.6	74.9	152
2× ↓ & ↑	50.2	6.4	32.3	56.8	74.6	152
4× ↓ & ↑	44.0	3.3	21.1	48.7	72.1	152
1333×800	52.4	8.5	37.4	58.6	74.9	152
2000×1200	59.8	22.2	48.6	64.4	77.4	182
2666×1600	62.7	30.7	53.3	66.4	77.8	287

### E. Experiments on the TJU-DHD-campus

In this subsection, some experiments on the TJU-DHD-campus are conducted. Table VI shows the average precisions (AP) of these four detectors. The anchor-free FCOS is 0.9% better than the anchor-based RetinaNet. Compared with the one-stage method, the two-stage method performs better. For example, FPN is 4.0% better and 3.1% better than FCOS and RetinaNet. Compared with FPN, cascade structure further improves detection performance by 2.7%. By comparing Table IV and VI, it can be seen that the same detector has a better performance on the TJU-DHD-traffic. The reason can be explained as follows. Compared with TJU-DHD-traffic, TJU-DHD-campus has a rich variance in objects (234,455 vs 27,650), seasons (4 vs 1), and image resolutions. Thus, TJU-DHD-campus is relatively difficult compared to TJU-DHD-traffic.

Table VII further shows average precisions (AP) on pedestrian category and rider category. It can be seen that: (1) On TJU-DHD-campus, the detectors also achieve a lower AP on the category of rider. For example, AP of RetinaNet on pedestrian detection is 4.1% higher than that on rider detection. (2) Cascade structure has a stable improvement on both pedestrian detection and rider detection. For example, Cascade R-CNN is 2.8% better and 2.7% better than FPN on pedestrian detection and rider detection, respectively.

To show the advantage of our high-resolution object dataset, we conduct two different experiments based on the detector

TABLE IX  
MISS RATES (MR) ON THE TJU-DHD-PEDESTRIAN. AS A REFERENCE, MISS RATE ON CITYPERSONS IS ALSO SHOWN. IN EACH DATASET, RESULTS ON DIFFERENT SUBSETS ARE SHOWN.

Method	set	R	RS	HO	R+HO	A
(a) RetinaNet [41]	TJU-Ped-campus	34.73	82.99	71.31	42.26	44.34
(b) FCOS [59]	TJU-Ped-campus	31.89	81.28	69.04	39.38	41.62
(c) FPN [40]	TJU-Ped-campus	27.92	73.14	67.52	35.67	38.08
(d) RetinaNet [41]	TJU-Ped-traffic	23.89	37.92	61.60	28.45	41.40
(e) FCOS [59]	TJU-Ped-traffic	24.35	37.40	63.73	28.86	40.02
(f) FPN [40]	TJU-Ped-traffic	22.30	35.19	60.30	26.71	37.78
(g) RetinaNet [41]	Citypersons	15.99	28.54	49.65	32.36	43.86
(h) FCOS [59]	Citypersons	18.29	27.70	52.42	34.73	44.39
(i) FPN [40]	Citypersons	14.25	26.67	49.26	29.71	38.79

FPN in Table VIII. One is first downsampling the images and then upsampling them to their original resolutions. We adopt 2× and 4× resampling (denoting as the downsampling/upsampling procedure for simplicity), which leads to 2.2% and 8.4% drop in AP, respectively. Moreover, the resolution plays a more important role in small-scale object detection. The AP<sub>s</sub> drops by 2.2% with 2× resampling and 16.3% with 4× resampling. It demonstrates that high-resolution image quality is useful for improving object detection performance, especially for small-scale object detection performance.

The second experiment is using different input sizes (*i.e.*, 1,333×800 pixels, 2,000×1,200 pixels, and 2,666×1,600 pixels). With the increment of input size, the detection accuracy becomes higher. When the input size ranges from 1,333×800 pixels to 2,000×1,200 pixels, the AP has 7.4% improvement. Moreover, the improvement mainly comes from tiny and small objects. When the input size ranges from 1,333×800 pixels to 2,000×1,200 pixels, AP<sub>t</sub> and AP<sub>s</sub> have 13.7% and 11.2% improvements, while AP<sub>m</sub> and AP<sub>l</sub> have 5.8% and 2.5% improvements. When using the high-resolution 2,666×1,600 pixels, it totally has a gain of 10.3% on AP and a gain of 22.2% on AP<sub>t</sub>. It means that higher resolution is very useful to detect small-scale objects. However, detecting objects in higher-resolution images needs much more computational cost. Thus, designing efficient detectors for high-resolution images becomes necessary in the future.

### F. Experiments on the TJU-DHD-pedestrian

In this subsection, some experiments based on RetinaNet [41], FCOS [59], and FPN [40] are conducted on the TJU-DHD-pedestrian. Similar to the Citypersons dataset [75], miss rates under the different subsets are shown. These subsets are the reasonable set (R), the reasonable small set (RS), the heavy

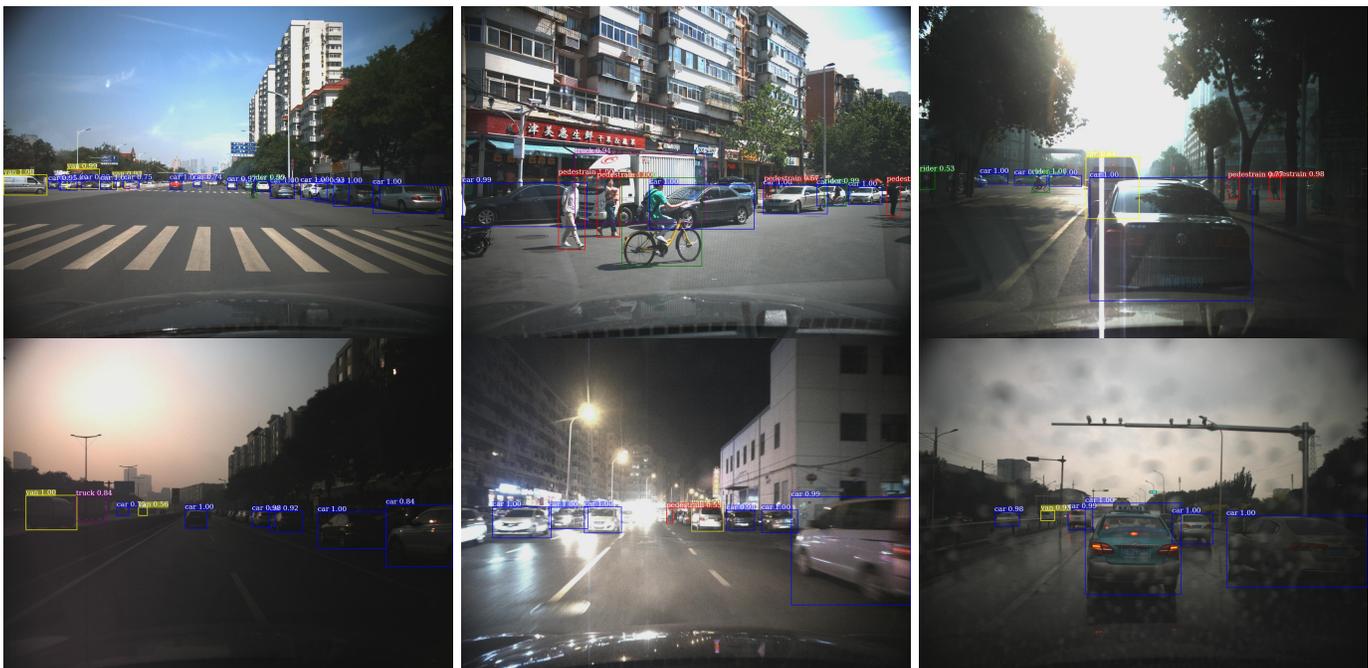


Fig. 7. Qualitative results of Cascade R-CNN on the TJU-DHD-traffic. The images under different illuminations and different weathers are chosen.

TABLE X  
MISS RATES (MR) OF FPN BY USING CROSS-SCENE EVALUATION. TJU-PED-CAMPUS AND PED-TRAFFIC IN TJU-DHD-PEDESTRIAN, CITYPERSONS, EURO-CITY ARE FOUR DIFFERENT DATASETS.

train set	test set	<b>R</b>	<b>R+HO</b>
(a) TJU-Ped-campus	TJU-Ped-campus	27.92	35.67
	TJU-Ped-traffic	30.62	33.89
	Citypersons	16.15	33.46
	Eurocity	14.10	26.73
	mean	22.20	32.43
(b) TJU-Ped-traffic	TJU-Ped-campus	42.08	50.55
	TJU-Ped-traffic	22.30	26.71
	Citypersons	23.79	43.09
	Eurocity	27.37	44.14
	mean	28.89	41.12
(c) Citypersons	TJU-Ped-campus	47.70	54.77
	TJU-Ped-traffic	46.27	49.40
	Citypersons	14.25	29.71
	Eurocity	21.99	35.58
	mean	32.55	42.37
(d) Eurocity	TJU-Ped-campus	41.86	48.52
	TJU-Ped-traffic	43.27	46.12
	Citypersons	16.16	31.51
	Eurocity	6.81	15.83
	mean	27.78	35.50

occlusion set (**HO**), the reasonable set and the heavy occlusion set (**R+HO**), and the all set (**A**).

Table IX shows miss rates of the same-scene evaluation on the TJU-DHD-pedestrian. The detector is trained and tested on the same scene dataset (*i.e.*, either TJU-Ped-campus or TJU-Ped-traffic). As a reference, miss rates on the Citypersons are also shown. It can be seen as follows: (1) Similar to object detection, the two-stage detector FPN outperforms the single-stage detector RetinaNet and the anchor-free detector FCOS on

pedestrian detection. For example, FPN outperforms RetinaNet and FCOS by 6.26% and 3.54% on the **A** set of TJU-Ped-campus. (2) TJU-Ped-campus and TJU-Ped-traffic are more challenging than Citypersons. For example, miss rates of FPN are 27.92% and 22.30% on the **R** set of TJU-Ped-campus and TJU-Ped-traffic, while the miss rate of FPN is 14.25% on the **R** set of Citypersons. Namely, miss rates of FPN on the **R** set of TJU-Ped-campus and TJU-Ped-traffic are 13.67% and 8.05% higher than that on the **R** set of Citypersons. Thus, there is a large space to improve the performance on our newly built pedestrian dataset.

To show the diversity of the built TJU-DHD-pedestrian quantitatively, Table X further gives a cross-scene evaluation on the TJU-Ped-campus, TJU-Ped-traffic, Citypersons [75], and EuroCity persons [1] based on two-stage FPN. The miss rate on each dataset and the mean miss rate over the four datasets are both given. If one dataset has a better diversity than another dataset, we think that it will have a lower mean miss rate when performing the cross-scene evaluation. It can be seen as follows: (1) The detector trained on TJU-Ped-campus has a best diversity, which achieves the lowest mean miss rate. For example, the detector trained on TJU-Ped-campus achieves 22.20% mean miss rate on the **R**, while that trained on Citypersons achieves 32.55% mean miss rate on the **R**. Namely, the detector trained on TJU-Ped-campus outperforms that trained on Citypersons by 10.35%. Similarly, TJU-Ped-campus outperforms EuroCity by 5.58%. It is demonstrated that our TJU-DHD-pedestrian has a richer diversity compared to EuroCity. (2) Our TJU-Ped-campus achieves a lower mean miss rate than our TJU-Ped-traffic. The reason can be explained as follows: Compared with TJU-Ped-traffic, TJU-Ped-campus has a rich variance in objects (234,455 vs 27,650) and seasons (4 vs 1). (3) Though TJU-

TABLE XI  
MISS RATES (MR) ON THE CITYPERSONS BY FINE-TUNING FPN ON OUR TJU-DHD-PEDESTRIAN.

Method	R	RS	HO	R+HO	A
(a) Citypersons [75]	14.25	26.67	49.26	29.71	38.79
(b) only TJU-Ped-traffic $\rightarrow$ Citypersons	13.23	22.81	48.93	29.01	37.38
(c) only TJU-Ped-campus $\rightarrow$ Citypersons	10.68	18.99	41.15	24.40	32.44
(d) TJU-DHD-pedestrian $\rightarrow$ Citypersons	10.13	17.07	40.40	23.88	31.89



Fig. 8. Qualitative results of Cascade R-CNN on TJU-DHD-campus. The images under different illuminations and different seasons are chosen.



Fig. 9. Some failure cases on the TJU-DHD-traffic and TJU-DHD-campus. The top row shows the false objects, and the bottom row shows the missed objects.

Ped-traffic and EuroCity are collected from the traffic scenes, they have a large domain gap. For example, the detectors trained on these two datasets have 20.97% difference when testing on TJU-Ped-traffic, 20.56% difference when testing on EuroCity. It means that the traffic scenes in China and Europe have a large difference. Thus, these two datasets are complementary.

Finally, we fine-tune the two-stage detector FPN on the

recently built Citypersons by using our TJU-DHD-pedestrian in Table XI. Namely, the detector is firstly trained based on our TJU-DHD-pedestrian and secondly fine-tuned on the Citypersons dataset. The initial miss rate on the **R** set of Citypersons is 14.25%. If only using TJU-Ped-traffic, miss rate on the **R** set drops from 14.25% to 13.23%. If only using TJU-Ped-campus, miss rate on the **R** set drops from 14.25% to 10.68%. If using both TJU-Ped-traffic and TJU-Ped-campus,

miss rate on the **R** set drops from 14.25% to 10.13%. The total improvement using our TJU-DHD-pedestrian is 4.12% on the **R** set. Namely, our TJU-DHD-pedestrian can help improve the detection performance on the Citypersons dataset.

### G. Visualizations

In this subsection, some visualizations on both TJU-DHD-traffic and TJU-DHD-campus are given. Because pedestrian detection is a special case of object detection in our dataset, the visualizations of pedestrian detection results on the TJU-DHD-pedestrian are not further given. Fig. 7 and Fig. 8 respectively show detection results of Cascade R-CNN [5] on the TJU-DHD-traffic and TJU-DHD-campus. To better see the performance of Cascade R-CNN on different conditions, the images under different illumination variance, different weathers, and different seasons are chosen for visualizations. It can be seen that Cascade R-CNN achieves a good result under these variances in some degree.

Meanwhile, some failure cases, including the false positives and the false negatives, of Cascade R-CNN on the TJU-DHD-traffic and TJU-DHD-campus are further given in Fig. 9. The top row shows examples of the false objects. It can be seen that some object-like things (e.g., air-conditioners and street signs), object-like background (e.g., potted plants), and some parts of objects are easily recognized as the detected objects. The bottom row gives examples of missed objects. It can be seen that the occluded objects, the small-scale objects, and low-illumination objects are usually missed by the detector. These problems remain unsolved and might be the key for further research in object detection.

### H. Discussion

*Challenges* Based on the above qualitative and quantitative experiments, we give a summary about the main challenges of object detection and pedestrian detection as follows: (1) Small-scale object detection and occluded object detection are still the bottlenecks of object detection. On the one hand, small-scale objects and occluded objects both have limited useful information. Small-scale objects have a noisy and blurred appearance and contain limited useful information, while occluded objects lose some local part information. On the other hand, small-scale objects and occluded objects make the intra-class variance large. (2) Most of the object detectors have a poor generality ability. For example, the detector trained on EuroCity persons has a poor performance on our TJU-Ped-campus (see Table X). However, to meet the actual application requirements, the object detector needs to have good robustness for domain adaption.

*Future Directions* Based on the challenges mentioned above, some future directions are summarized as follows: (1) *Small-scale object detection*. One direction is that designing better feature pyramid and image pyramid structures. Another direction is that designing efficient detectors for high-resolution object detection. As a result, it does not add much computational cost when using high-resolution images. (2) *Occluded object detection*. One direction is that making full use of visible regions to suppress the negative impact of occluded

regions. Another direction is that exploiting how to retain correct detection in the crowded object detection. (3) *Cross-scene object detection*. Traffic scene and campus scene are two common scenes, which have important application value. How to make the detector have a good generation ability is an important topic in the future.

## V. CONCLUSION

In this paper, a new diverse high-resolution dataset for object detection in two typical scenes (traffic scene and campus scene) was built. The experiments based on four typical detectors are conducted to give the baseline performance on this dataset. Meanwhile, based on this dataset, a new large-scale pedestrian dataset is also built. The built datasets, TJU-DHD-traffic, TJU-DHD-campus, and TJU-DHD-pedestrian, can hopefully be one of the fundamental benchmarks for object detection and pedestrian detection.

## REFERENCES

- [1] M. Braun, S. Krebs, F. Flohr, and D. M. Gavrilu, "The EuroCity persons dataset: A novel benchmark for object detection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 41, no. 8, pp. 1844-1861, 2019. **2, 5, 6, 9**
- [2] G. Brazil, X. Yin, and X. Liu, "Illuminating pedestrians via simultaneous detection and segmentation," *Proc. IEEE International Conference on Computer Vision*, 2017. **3**
- [3] G. Brazil and X. Liu, "Pedestrian detection With autoregressive network phases," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. **3**
- [4] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," *Proc. European Conference on Computer Vision*, 2016. **3**
- [5] Z. Cai and N. Vasconcelos, "Cascade r-cnn: Delving into high quality object detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018. **2, 6, 7, 8, 11**
- [6] J. Cao, Y. Pang, and X. Li, "Triply supervised decoder networks for joint detection and segmentation," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. **3**
- [7] J. Cao, Y. Pang, and X. Li, "Pedestrian detection inspired by appearance constancy and shape symmetry," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016. **3, 5**
- [8] J. Cao, Y. Pang, and X. Li, "Learning multilayer channel features for pedestrian detection," *IEEE Transactions on Image Processing*, 2016. **3**
- [9] J. Cao, Y. Pang, J. Han, and X. Li, "Hierarchical shot detector," *Proc. IEEE International Conference on Computer Vision*, 2019. **3**
- [10] J. Cao, Y. Pang, S. Zhao, and X. Li, "High-level semantic networks for multi-scale object detection," *IEEE Trans. on Circuits and Systems for Video Technology*, 2019. **2, 3**
- [11] C. Chi, S. Zhang, J. Xing, Z. Lei, S. Z. Li, and X. Zou, "PedHunter: Occlusion robust pedestrian detector in crowded scenes," *Proc. AAAI Conference on Artificial Intelligence*, 2020. **3**
- [12] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The Cityscapes dataset for semantic urban scene understanding," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016. **2, 6**
- [13] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2005. **2**
- [14] P. Dollár, Z. Tu, P. Perona, and S. Belongie, "Integral channel features," *Proc. British Machine Vision Conference*, 2009. **2, 3**
- [15] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: An evaluation of the state of the art," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 34, no. 4, pp. 743-761, 2012. **1, 2, 4, 5, 6**
- [16] A. Ess, B. Leibe, and L. Van Gool, "Depth and appearance for mobile scene analysis," *Proc. IEEE International Conference on Computer Vision*, 2007. **2**
- [17] M. Enzweiler and D. M. Gavrilu, "Monocular pedestrian detection: Survey and experiments," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 31, no. 12, pp. 2179-2195, 2009. **2**

- [18] M. Everingham, L. V. Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International Journal of Computer Vision*, vol. 88, no. 2, pp. 303–338, 2010. 1, 2, 6
- [19] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "CenterNet: Keypoint triplets for object detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [20] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 1, 2, 5, 6
- [21] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2014. 2, 3, 6
- [22] R. Girshick, "Fast r-cnn," *Proc. IEEE International Conference on Computer Vision*, 2015. 2, 3, 6
- [23] A. Gupta, P. Dollar, and R. Girshick, "LVIS: A dataset for large vocabulary instance segmentation," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2
- [24] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2, 3, 7
- [25] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Proc. Advances in Neural Information Processing Systems*, 2015. 2, 3, 6, 7
- [26] T. Kong, F. Sun, A. Yao, H. Liu, M. Lu, and Y. Chen, "Ron: Reverse connection with objectness prior networks for object detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [27] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Proc. Advances in Neural Information Processing Systems*, 2012. 1, 2
- [28] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *Proc. European Conference on Computer Vision*, 2014. 3
- [29] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 1, 2, 7
- [30] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," *Proc. IEEE International Conference on Computer Vision*, 2017. 2, 3
- [31] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2
- [32] X. Huang, P. Wang, X. Cheng, D. Zhou, Q. Geng, and R. Yang, "The ApolloScape open dataset for autonomous driving and its application," *IEEE Trans. Pattern Analysis and Machine Intelligence*, 2019. 1, 2
- [33] X. Huang, Z. Ge, Z. Jie, and O. Yoshie, "NMS by representative region: Towards crowded pedestrian detection by proposal pairing," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 3
- [34] V. Jain and E. Learned-Miller, "FDDB: A benchmark for face detection in unconstrained settings," *Technical report, University of Massachusetts*, 2010. 2
- [35] B. Jiang, R. Luo, J. Mao, T. Xiao, and Y. Jiang, "Acquisition of localization confidence for accurate object detection," *Proc. European Conference on Computer Vision*, 2018. 3
- [36] A. Kuznetsova, H. Rom, N. Alldrin, J. Uijlings, I. Krasin, J. Pont-Tuset, S. Kamali, S. Popov, M. Mallocci, T. Duerig, and V. Ferrari, "The Open Images Dataset V4: Unified image classification, object detection, and visual relationship detection at scale," *arXiv:1811.00982*, 2018. 2
- [37] Hei Law and Jia Deng, "CornerNet: Detecting objects as paired keypoints," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [38] Y. Li, Y. Chen, N. Wang, and Z. Zhang, "Scale-Aware Trident Networks for Object Detection," *Proc. IEEE International Conference on Computer Vision*, 2018. 3
- [39] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and L. Zitnick, "Microsoft coco: Common objects in context," *Proc. European Conference on Computer Vision*, 2014. 1, 2, 6
- [40] T. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 3, 6, 7, 8
- [41] T. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *Proc. IEEE International Conference on Computer Vision*, 2017. 2, 3, 6, 7, 8
- [42] C. Lin, J. Lu, G. Wang, and J. Zhou, "Graininess-Aware Deep Feature Learning for Pedestrian Detection," *Proc. European Conference on Computer Vision*, 2018. 3
- [43] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," *Proc. European Conference on Computer Vision*, 2016. 2, 3, 7
- [44] W. Liu, S. Liao, W. Ren, W. Hu, and Y. Yu, "High-level semantic feature detection: A new perspective for pedestrian detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [45] S. Liu, D. Huang, and Y. Wang, "Adaptive NMS: Refining pedestrian detection in a crowd," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 3
- [46] J. Mao, T. Xiao, Y. Jiang, and Z. Cao, "What can help pedestrian detection?," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3
- [47] M. Najibi, B. Singh, and L. S. Davis, "AutoFocus: Efficient multi-scale inference," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [48] G. Neuhof, T. Ollmann, S. Rota Bulo, and P. Kotschieder, "The Mapillary vistas dataset for semantic understanding of street scenes," *Proc. IEEE International Conference on Computer Vision*, 2017. 2
- [49] J. Nie, R. M. Anwer, H. Cholakkal, F. Shahbaz Khan, Y. Pang, and L. Shao, "Enriched Feature Guided Refinement Network for Object Detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 2
- [50] T. Ojala, M. Pietikainen, and T. Maenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002. 2
- [51] Y. Pang, J. Xie, M. H. Khan, R. M. Anwer, F. Shahbaz Khan, and L. Shao, "Mask-guided attention network for occluded pedestrian detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [52] Y. Pang, J. Cao, J. Wang, and J. Han, "JCS-Net: Joint classification and super-resolution network for small-scale pedestrian detection in surveillance images," *IEEE Trans. Information Forensics and Security*, vol. 14, no. 12, pp. 3322–3331, 2019. 3
- [53] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "Imagenet large scale visual recognition challenge," *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015. 1
- [54] S. Shao, Z. Zhao, B. Li, T. Xiao, G. Yu, X. Zhang, and J. Sun, "CrowdHuman: A benchmark for detecting human in a crowd," *arXiv:1805.00123*, 2018. 2, 5, 6
- [55] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *Proc. International Conference on Learning Representations*, 2015. 1, 2
- [56] B. Singh and L. S. Davis, "An analysis of scale invariance in object detection – snip," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [57] B. Singh, M. Najibi, and L. S. Davis, "Sniper: Efficient multi-scale training," *Proc. Advances in Neural Information Processing Systems*, 2018. 3
- [58] T. Song, L. Sun, D. Xie, H. Sun, and S. Pu, "Small-scale pedestrian detection based on topological line localization and temporal feature aggregation," *Proc. European Conference on Computer Vision*, 2018. 3
- [59] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 2, 3, 6, 7, 8
- [60] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders, "Selective search for object recognition," *International Journal of Computer Vision*, vol. 104, no. 2, pp. 154–171, 2013. 3
- [61] P. Viola and M. Jones, "Robust real-time face detection," *International Journal of Computer Vision*, vol. 57, no. 2, pp. 137–154, 2004. 2
- [62] X. Wang, T. Han, and S. Yan, "An HOG-LBP human detector with partial occlusion handling," *Proc. IEEE International Conference on Computer Vision*, 2009. 2
- [63] X. Wang, T. Xiao, Y. Jiang, S. Shao, J. Sun, and C. Shen, "Repulsion loss: Detecting pedestrians in a crowd," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [64] T. Wang, R. M. Anwer, H. Cholakkal, F. S. Khan, Y. Pang, and L. Shao, "Learning rich features at high-speed for single-shot object detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [65] J. Wu, C. Zhou, M. Yang, Q. Zhang, Y. Li, and J. Yuan, "Temporal-context enhanced detection of heavily occluded pedestrians," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 3
- [66] J. Wu, C. Zhou, Q. Zhang, M. Yang, and J. Yuan, "Self-Mimic Learning for Small-scale Pedestrian Detection," *Proc. ACM International Conference on Multimedia*, 2020. 3

- [67] T. Xiao, S. Li, B. Wang, L. Lin, and X. Wang, Joint Detection and Identification Feature Learning for Person Search *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [68] M. Ye, J. Li, A. J. Ma, L. Zheng, and P. Yuen, "Dynamic graph co-matching for unsupervised video-based person re-identification," *IEEE Trans. Image Processing*, vol. 28, no. 6, pp. 2976-2990, 2019. 2
- [69] M. Ye, X. Lan, Z. Wang, and P. Yuen, "Bi-directional center-constrained top-ranking for visible thermal person re-identification," *IEEE Trans. Information Forensics and Security*, 10.1109/TIFS.2019.2921454, 2019. 2
- [70] S. Yang, P. Luo, C. Change Loy, and X. Tang, "WIDER FACE: A face detection benchmark," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2016. 2
- [71] J. Yan, X. Zhang, Z. Lei, and S. Z. Li, "Face detection by structural models," *Image and Vision Computing*, 2014. 2
- [72] Z. Yang, S. Liu, H. Hu, L. Wang, and S. Lin, "RepPoints: Point set representation for object detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [73] F. Yu, W. Xian, Y. Chen, F. Liu, M. Liao, V. Madhavan, and T. Darrell, "BDD100K: A diverse driving video database with scalable annotation tooling," *arXiv:1805.04687*, 2019. 1, 2
- [74] S. Zhang, R. Benenson, and B. Schiele, "Filtered Channel Features for Pedestrian Detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 3, 5
- [75] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 1, 2, 4, 5, 6, 8, 9, 10
- [76] S. Zhang, J. Yang, and B. Schiele, "Occluded Pedestrian Detection Through Guided Attention in CNNs," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 3
- [77] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Single-shot refinement neural network for object detection," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [78] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Occlusion-aware R-CNN: Detecting pedestrians in a crowd," *Proc. European Conference on Computer Vision*, 2018. 3
- [79] Z. Zhang, S. Qiao, C. Xie, W. Shen, B. Wang, and A. L. Yuille, "Single-shot object detection with enriched semantics," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 3
- [80] Q. Zhao, T. Sheng, Y. Wang, Z. Tang, Y. Chen, L. Cai, and H. Ling, "M2det: A single-shot object detector based on multi-level feature pyramid network," *Proc. AAAI Conference on Artificial Intelligence*, 2018. 3
- [81] X. Zhou, D. Wang, and P. Krahenbühl, "Objects as points," *arXiv:1904.07850*, 2019. 3
- [82] C. Zhou, M. Yang, and J. Yuan, "Discriminative feature transformation for occluded pedestrian detection," *Proc. IEEE International Conference on Computer Vision*, 2019. 3
- [83] C. Zhou and J. Yuan, "Bi-box regression for pedestrian detection and occlusion estimation," *Proc. European Conf. on Computer Vision*, 2018. 3
- [84] C. Zhou and J. Yuan, "Multi-label learning of part detectors for heavily occluded pedestrian detection," *Proc. IEEE International Conference on Computer Vision*, 2017. 3
- [85] X. Zhu and D. Ramanan, "Face detection, pose estimation, and landmark localization in the wild," *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2012. 2