

WiderPerson: A Diverse Dataset for Dense Pedestrian Detection in the Wild

Shifeng Zhang, Yiliang Xie, Jun Wan, *Member, IEEE*, Hansheng Xia, and Stan Z. Li, *Fellow, IEEE*, and Guodong Guo, *Senior, IEEE*

Abstract—Pedestrian detection has achieved significant progress with the availability of existing benchmark datasets. However, there is a gap in the *diversity* and *density* between real world requirements and current pedestrian detection benchmarks: 1) most of existing datasets are taken from a vehicle driving through the regular traffic scenario, usually leading to insufficient diversity; 2) crowd scenarios with highly occluded pedestrians are still under represented, resulting in low density. To narrow this gap and facilitate future pedestrian detection research, we introduce a large and diverse dataset named WiderPerson for dense pedestrian detection in the wild. This dataset involves five types of annotations in a wide range of scenarios, no longer limited to the traffic scenario. There are a total of 13,382 images with 399,786 annotations, *i.e.*, 29.87 annotations per image, which means this dataset contains dense pedestrians with various kinds of occlusions. Hence, pedestrians in the proposed dataset are extremely challenging due to large variations in the scenario and occlusion, which is suitable to evaluate pedestrian detectors in the wild. We introduce an improved Faster R-CNN and the vanilla RetinaNet to serve as baselines for the new pedestrian detection benchmark. Several experiments are conducted on previous datasets including Caltech-USA and CityPersons to analyze the generalization capabilities of the proposed dataset and we achieve state-of-the-art performances on these previous datasets without bells and whistles. Finally, we analyze common failure cases and find the classification ability of pedestrian detector needs to be improved to reduce false alarm and miss detection rates. The proposed dataset is available at <http://www.cbsr.ia.ac.cn/users/sfzhang/WiderPerson>.

Index Terms—Pedestrian detection, dataset, rich diversity, high density.

I. INTRODUCTION

PEDESTRIAN detection is a long-standing problem in computer vision and pattern recognition with extensive applications including security and surveillance, mobile robotics, autonomous driving, and crowd sourcing, to name a few. The accuracy of pedestrian detection systems has a direct impact on these tasks, hence the success of pedestrian

Shifeng Zhang, Jun Wan and Stan Z. Li are with the Center for Biometric Security Research (CBSR), National Laboratory of Pattern Recognition (NLPR), Institute of Automation Chinese Academy of Sciences (CASIA) and University of Chinese Academy of Sciences (UCAS), Beijing, China. Stan Z. Li is also with the Macau University of Science and Technology, Macau, China (e-mail: {shifeng.zhang, jun.wan, szli}@nlpr.ia.ac.cn).

Yiliang Xie is with the University of Southern California (USC), US (e-mail: microos316@gmail.com).

Hansheng Xia is with the College of Energy and Power Engineering, Nanjing University of Aeronautics and Astronautics (NUAA), Nanjing, China (e-mail: hanson_cha@163.com).

Guodong Guo is with the Institute of Deep Learning, Baidu Research and National Engineering Laboratory for Deep Learning Technology and Application (e-mail: guogudong01@baidu.com).



Fig. 1. The diversity and density of the newly introduced WiderPerson dataset. It can bridge the gap between real world requirements and pedestrian detection benchmarks. For visualization, we use bounding boxes of different colors for pedestrians (Cyan), riders (Red), partially-visible persons (Green), crowd (Yellow) and ignore regions (Blue).

detection is of crucial importance. Given an arbitrary image, the goal of pedestrian detection is to determine whether or not there are any pedestrians in the image, and if present, return the image location and extent of each pedestrian. While this appears as an effortless task for human, it is a very difficult task for computers. The challenges associated with pedestrian detection can be attributed to variations in pose, scale and occlusion, which need to be addressed while building pedestrian detection algorithms.

With the remarkable progress over the past few decades, pedestrian detection has been successfully applied in some practical application systems under restricted scenarios. The success of these systems can be attributed to two key steps: (1) advancements in the field of deep Convolutional Neural Network (CNN) which has had a direct impact on many computer vision tasks including pedestrian detection; (2) dataset

TABLE I

COMPARISON OF PEDESTRIAN DETECTION DATASETS (Training SUBSET ONLY). ‘-’ MEANS THIS TERM IS UNLIMITED AND CAN NOT BE COUNTED.

	Caltech-USA	KITTI	CityPersons	WiderPerson
# country	1	1	3	—
# city	1	1	18	—
# season	1	1	3	4
# images	42,782	3,712	2,975	8,000
# persons	13,674	2,322	19,654	236,073
# ignore regions	50,363	45	6,768	8,979
# person/image	0.32	0.63	6.61	29.51
# unique persons	1,273	< 2,322	19,654	236,073

collection efforts led by different researchers in the community. Furthermore, improvements in detection algorithms have almost always been followed by the publication of more challenging datasets and vice versa. Such synchronous advancement in both steps has led to an even more rapid progress in the field. In terms of pedestrian detection, publicly available benchmark datasets such as Caltech-USA [1], KITTI [2] and CityPersons [3] have contributed to spurring interest and progress in pedestrian detection research. Coupled with the development and blooming of deep learning, modern pedestrian detectors [4]–[8] have achieved remarkable performance.

Although performance has been significantly improved, it’s still difficult to assess for real world, since compared with crowd-counting datasets [9], [10] designed in the crowded condition, there is a gap in the diversity and density between current existing pedestrian detection benchmarks and real world requirements. On the one hand, most of existing datasets are collected via a vehicle-mounted camera through the regular traffic scenario. This fixed scenario significantly reduces the richness of the foreground and background, leading to low diversity. Specifically, only pedestrians and backgrounds on the road are taken into consideration while the other scenarios are severely under represented. Thus, diversity in pedestrian and background appearances is limited. On the other hand, crowd scenarios with highly occluded pedestrians are still under-represented. As shown in Table I, the Caltech-USA and KITTI datasets have less than one person per image, while the CityPersons dataset has ~ 7 persons per image. Even worse, protocols of these datasets allow annotators to ignore the regions with a large number of persons, since exhaustively annotating crowd regions is incredibly difficult and time consuming, resulting in low density and insufficient occlusions cases. To sum up, current pedestrian detection datasets typically contain a few thousand pedestrians with limited variations in diversity and density. These limitations have partially contributed to the failure of some algorithms in coping with heavy occlusion and atypical scenario. Therefore, more challenging datasets similar to the real world are needed to trigger progress and inspire novel ideas.

To move forward the field of pedestrian detection, we introduce a diverse and dense pedestrian detection dataset called WiderPerson. It consists of 13,382 images with 399,786 annotations, *i.e.*, 29.87 annotations per image, varying largely in scenario and occlusion, as shown in Fig. 1. Besides, the annotations have five fine-grained labels, *i.e.*, pedestrians,

riders, partially-visible persons, crowd, and ignore regions. These high quality annotations provide a rich diverse dataset and enable new experiments both for training better models, and as new test benchmark. We split the proposed WiderPerson dataset into three subsets (training, validation, and testing sets). Annotations of training and validation will be public, and an online benchmark will be set-up. We show an example of using the proposed WiderPerson dataset through proposing an improved Faster R-CNN [11], which consists of finer feature map, ignore region and tiny pedestrian handling, Region of Interest (RoI) feature enhancing and dynamic sample strategy to deal with large density and diversity variations. The cross-dataset generalization results of the proposed WiderPerson dataset show that it is an effective training source for pedestrian detection and we achieve state-of-the-art performance on existing Caltech-USA and CityPersons datasets.

For clarity, the main contributions of this work can be summarized as three-fold:

- We propose the WiderPerson dataset, which provides a large number of highly diverse and dense bounding box annotations for pedestrian detection.
- We build an improved Faster R-CNN to show an example of using WiderPerson, which consists of some improvements to deal with large density and diversity variations.
- We prove the generalization capabilities of detectors trained with the new dataset and achieve state-of-the-art performance on Caltech-USA and CityPersons datasets.

The rest of the paper is organized as follows. Section II reviews the related work. Description of the WiderPerson dataset is presented in Section III. Section IV introduces our proposed baseline detector and Section V shows the experimental results. Section VI concludes the paper.

II. RELATED WORK

A. Dataset

In the last decade, several datasets have been created for pedestrian detection training and evaluation. The GM-ATCI dataset [12] is collected using a vehicle-mounted standard automotive rear-view display camera for evaluating rear-view pedestrian detection. The INRIA dataset [13] is one of the most popular static pedestrian detection datasets. The USC dataset [14] consists of a number of fairly small pedestrian datasets taken largely from surveillance video. The ETH dataset [15] is captured from a stereo rig mounted on a stroller in the urban. The CVC-ADAS dataset [16] contains pedestrian videos acquired on-board, virtual-world pedestrians (with part annotations) and occluded pedestrians. The NICTA dataset [17] is a large scale urban dataset collected in multiple cities and countries, it has no motion and tracking information but significant number of unique pedestrians. The Daimler dataset [18] is captured in an urban setting and has tracking information and a large number of labelled bounding boxes. The TUD-Brussels dataset [19] contains image pairs recorded in a crowded urban setting with an onboard camera. These datasets represent early efforts to collect pedestrian datasets.

Although these early datasets have contributed to spurring interest and progress of pedestrian detection, however, as algorithm performance improves, they are replaced by the larger

and richer datasets. The Tsinghua-Daimler Cyclist (TDC) dataset [20] focuses on cyclists recorded from a vehicle-mounted stereo vision camera, containing a large number of cyclists varying widely in appearance, pose, scale, occlusion and viewpoint. In [21], a multi-spectral dataset for pedestrian detection is introduced, combining RGB and infrared modalities. The Caltech-USA [1] dataset consists of approximately 10 hours of 640×480 30Hz video taken from a vehicle driving through regular traffic in an urban environment, which has been extended by [22] with corrected annotations. The KITTI [2] dataset focuses autonomous driving and is collected via a standard station wagon with two high-resolution color and grayscale video cameras, around the mid-size city of Karlsruhe, in rural areas and on highways, up to 15 cars and 30 pedestrians are visible per image. The CityPersons [3] dataset is recorded by a car traversing 27 different cities and provides high quality bounding boxes with larger portions of occluded persons. The EuroCity Persons dataset [23] provides a large number of highly diverse, accurate and detailed annotations of pedestrians, cyclists and other riders in 31 cities of 12 European countries.

Despite the prevalence of these datasets, they all suffer a problem of low diversity. Most of existing datasets are collected via a vehicle-mounted camera through the regular traffic scenario. The diversity in pedestrian and background appearances is limited. Another weakness of both datasets is that the crowd scenarios are significantly under represented, resulting in insufficient occlusions cases. The paper aims at solving these two problems via proposing a diverse and dense pedestrian detection dataset, which can narrow the gap in the diversity and density between real world requirements and current pedestrian detection benchmarks to better evaluate detectors in the wild. Besides, the proposed dataset is also very useful for training a re-detector for dealing with tracking loss for pedestrian tracking [24], [25].

B. Method

Generic Object Detection. Early generic object detection methods rely on the sliding window paradigm based on the hand-crafted features and classifiers to find the objects of interest. In recent years, with the advent of deep Convolutional Neural Network (CNN), a new generation of more effective object detection methods based on CNN significantly improve the state-of-the-art performances, which can be roughly divided into two categories, *i.e.*, the one-stage approach and the two-stage approach. The one-stage approach [26], [27] directly predicts object class label and regresses object bounding box based on the pre-tiled anchor boxes using deep CNN. The main advantage of the one-stage approach is its high computational efficiency. In contrast to the one-stage approach, the two-stage approach [11], [28] always achieves top accuracy on several benchmarks, which first generates a pool of object proposals by a separated proposal generator, and then predicts the class label and accurate location and size of each proposal.

Pedestrian Detection. Even as one of the long-standing problems in computer vision field with an extensive literature, pedestrian detection still receives considerable interests with

a wide range of applications. A common paradigm [29]–[31] to deal with this problem is to train a pedestrian detector that exhaustively operates on the sub-images across all locations and scales. Dalal and Triggs [13] design the Histograms of Oriented Gradient (HOG) descriptors and Support Vector Machine (SVM) classifier for human detection. Dollár *et al.* [32] demonstrate that using features from multiple channels can greatly improve the performance. Zhang *et al.* [33] provide a systematic analysis for the filtered channel features, and find that with the proper filter bank, filtered channel features can reach top detection quality. Paisitkriangkrai *et al.* [34] design a new feature built on low-level features and spatial pooling, and directly optimize the partial area under the Receiver Operating Characteristic (ROC) curve for better performance.

Recently, CNN-based detectors [35]–[38] have become a predominating trend in the field of pedestrian detection. Sermanet *et al.* [35] present an unsupervised method using the convolutional sparse coding to pre-train CNN for pedestrian detection. In [39], a complexity-aware cascaded detector is proposed for an optimal trade-off between accuracy and speed. Angelova *et al.* [40] combine the ideas of fast cascade and a deep network to detect pedestrian. Yang *et al.* [41] use scale-dependent pooling and layer-wise cascaded rejection classifiers to detect objects efficiently. Zhang *et al.* [42] present an effective pipeline for pedestrian detection via extracting self-learned features from the Region Proposal Network (RPN) [11] followed by a boosted decision forest. Cai *et al.* [43] propose an architecture which uses different levels of features to detect persons at various scales. Mao *et al.* [44] present a multi-task network architecture to jointly learn pedestrian detection with the given extra features. Li *et al.* [7] use multiple built-in sub-networks to adaptively detect pedestrians across scales. Brazil *et al.* [38] exploit weakly annotated bounding boxes via a segmentation infusion network to achieve considerable performance gains.

Occlusion is one of the most significant challenges in compute vision, especially for pedestrian detection, which increases the difficulty in pedestrian localization. Several methods [45]–[49] use part-based model to describe the pedestrian in occlusion handling, which learn a series of part detectors and design some mechanisms to fuse the part detection results to localize partially occluded pedestrians. Besides the part-based model, Leibe *et al.* [50] propose an implicit shape model to generate a set of pedestrian hypotheses that are further refined to obtain the visible regions. Wang *et al.* [51] divide the template of pedestrian into a set of blocks and conduct occlusion reasoning by estimating the visibility status of each block. Ouyang *et al.* [52] exploit multi-pedestrian detectors to aid single-pedestrian detectors to handle partial occlusions, especially when the pedestrians gather together and occlude each other in real-world scenarios. In [53], a set of occlusion patterns of pedestrians are discovered to learn a mixture of occlusion-specific detectors. Zhou *et al.* [54] propose to jointly learn part detectors to exploit part correlations and reduce the computational cost. Wang *et al.* [5] introduce a new bounding box regression loss to detect pedestrians in crowd scenarios.

III. PROPOSED WIDERPERSON DATASET

In this section, we present our WiderPerson dataset from aspects of collection process, annotation tool, annotation method, various statistical information and benchmarking.

A. Data Collection

For the diversity of our dataset, we crawl images from multiple image search engines ranging from Google, Bing, and Baidu. Combined with specially-designed keywords, one of the prominent advantages of using different image search engines together is the collected images possess diverse features in cities, events, and scenarios. We design more than 50 keywords (*e.g.*, pedestrian, cyclist, walking, running, marathon, square dance and group photo) during the crawling process and obtain $\sim 50,000$ images as our candidate images. To prevent the duplication of images, we leverage a simple but powerful mechanism, the pHash [55], along with the union find, for the removal of the repetitions. Moreover, images with sparse distribution of people are filtered out to keep the difficulties of our dataset. Finally, we have 13,382 images remained, and they are randomly split into training, validation and testing subsets with 8,000, 1,000 and 4,382 images, respectively.

B. Annotation Tool

We design a new annotation tool whose Graphical User Interface (GUI) is illustrated in Fig. 2. It is written using JavaScript and built with a very responsive design. The list of images that need to be marked is displayed on the upper right side. For the selected image to be annotated, the tool displays five kinds of annotation examples on the left side to help annotators to mark. These five different types of annotations are labelled with different colors to better distinguish. All the labelled annotations are shown on the lower right side and annotators can select any labelled annotation to display and correct. To complete the annotation, the user shall next adjust the position of the bounding boxes. For this purpose, the keyboard arrows shall be used. More precisely, the left, right, up and down keys should be used in order to shift the annotations on the image. To help annotators, there are two buttons (display and hide labels) at the top used to make labelled annotations optionally visible while annotating. In addition, a zooming feature at the top can be used to zooming-in the corresponding image and annotations. This is implemented in order to make it easier for the annotators to obtain more precise locations of the annotations. After getting used to the annotation process, annotators become more and more precise on these steps, which significantly reduces the time required to annotate as fewer adjustments are required. Besides, to ensure that profit is not prioritized over the accuracy and the precision of the annotations, we are highly involved in the process and all annotators must pass the strict annotation testing.

C. Image Annotation

Our Annotations are finely classified into five categories: pedestrians, riders, partially-visible persons, crowd and ignore regions. The annotation process contains the two steps:

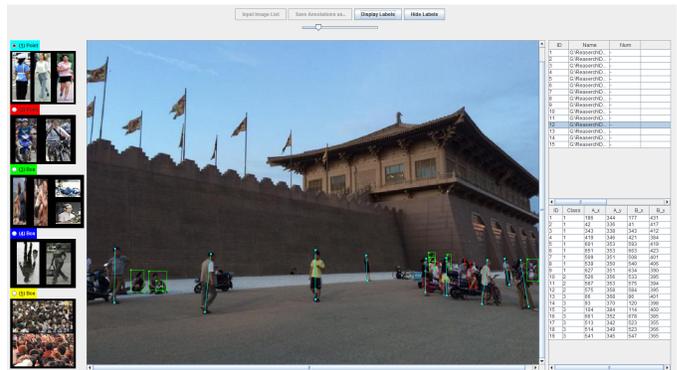


Fig. 2. Graphical User Interface (GUI) of our annotation tool.



Fig. 3. Illustration of bounding box annotations for pedestrians and riders. For each target, the top of the head and middle of the feet is drawn by the annotator. An aligned bounding box is automatically generated using the fixed aspect ratio (0.41).

1. Annotators are asked to thoroughly search across the whole image for individuals, and annotate them for using the similar protocol from [3]. For pedestrians and riders (shown in Fig. 3(a)), we generate a bounding box by drawing a line across one's head and the middle point between feet, as shown in Fig. 3(b). A bounding box aligned to the center of the line is then generated with an aspect ratio of 0.41 (defined as w/h), as shown in Fig. 3(c). For partially-visible persons, including individuals that are heavily-occluded or with unusual poses and viewpoints, we mark them using bounding boxes with unconstrained aspect ratios. The crowd in our dataset plays another critical role contributing to the variances and difficulties. Similar to the partially-visible persons, we also annotate a group of people using a tightly bounded rectangle. Finally, we annotate regions containing fake human, for instance, human on the posters, reflections, mannequin and statues, etc.
2. After the above-mentioned annotating process, to ensure the quality of the labels, we perform three-fold cross-validation to check the annotations strictly. Each image is intuitively marked as either correct or erroneous by three

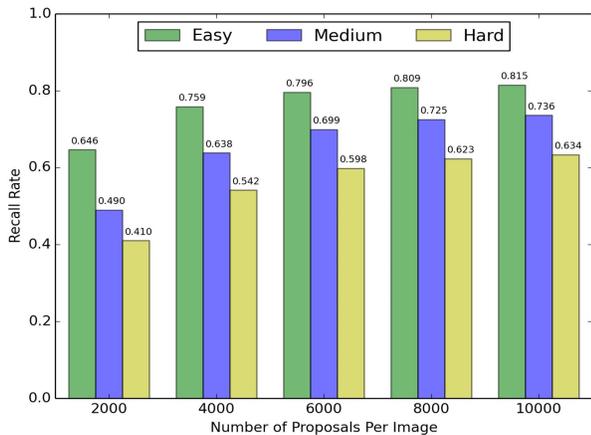


Fig. 4. Recall rate with different number of proposals. Proposals are generated by using Edgebox [57]. Lower recall rate implies higher difficulty. We show histograms of detection rate over the number of proposal for different subsets.

TABLE II
STATISTICS OF ANNOTATIONS ON WIDERPERSON DATASET.

	Training	Validation	Testing	Sum
# images	8,000	1,000	4,382	13,382
# persons	236,073	27,762	122,518	386,353
# ignore regions	8,979	661	3,793	13,433
# person/images	29.51	27.76	27.96	28.87

different annotators, and if it marked as erroneous by more than half of the annotators, it would be re-annotated until it passes the check. Fig. 1 shows some exemplary final annotations.

D. Dataset Statistic

Capacity. The number of bounding box annotations provided by our WiderPerson dataset is shown in table II, which illustrates the capacity of WiderPerson dataset. In a total of 13,382 images, there are $\sim 386k$ person and $\sim 13k$ ignore region annotations in the WiderPerson dataset. The number of annotations is more than $10\times$ boosted compared with previous challenging pedestrian detection dataset like CityPersons. The total number of persons is also noticeably larger than the others. We randomly select 8000/1000/4382 images as training/validation/testing subsets. Following the principle in WIDER FACE [56], we define three levels of difficulty: ‘Easy’ (≥ 100 pixels), ‘Medium’ (≥ 50 pixels), ‘Hard’ (≥ 20 pixels) according to the physical height of ground-truth bounding boxes. As shown in the Fig. 4, we utilize EdgeBox [57] to evaluate their detection rates with different number of proposals. The average recall rates for these three levels are 81.5%, 73.6% and 63.4% with 10,000 proposal per image.

Scale. To analyse the scale characteristic across different datasets, we use the probability density function (PDF) to specify the probability of scale falling within a particular range of values, which can specify the distribution of scales. To this end, we group the persons by their image size (height in pixels) into some scale bins. As can be observed from Fig. 5, Caltech-USA and CityPersons have a limited scale distribution, most

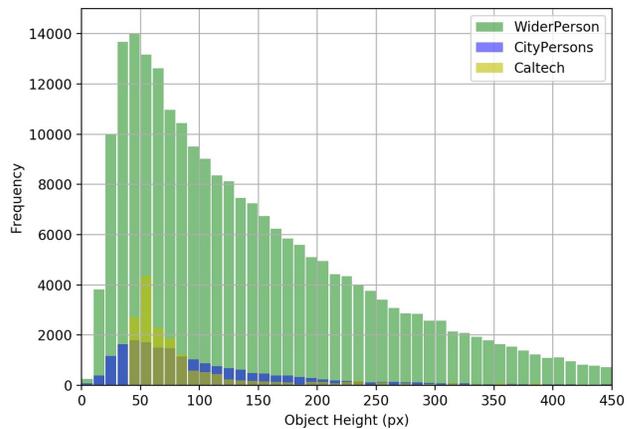


Fig. 5. Scale distribution of different dataset. We use the probability density function (PDF) to specify the probability of scale falling within a particular range of values.

TABLE III
DENSITY COMPARISON BETWEEN WIDELY USED PEDESTRIAN DETECTION DATASETS. IT DEMONSTRATES THE NUMBER AND PROPORTION OF IMAGES THAT CONTAIN \geq # PERSONS IN DIFFERENT DATASETS.

# persons	Caltech-USA	CityPersons	WiderPerson
≥ 1	7,839 18.3%	2,482 83.4%	8,000 100.0%
≥ 2	3,257 7.6%	2,082 70.0%	7,999 100.0%
≥ 3	1,265 3.0%	1,741 58.5%	7,998 100.0%
≥ 5	282 0.7%	1,225 41.2%	7,994 99.9%
≥ 10	36 0.1%	610 20.5%	7,924 99.1%
≥ 20	0 0.0%	227 7.6%	5,145 64.3%
≥ 30	0 0.0%	94 3.2%	2,564 32.1%

of their annotations are between 30~100 pixels in height. In contrast, our WiderPerson dataset covers a much wider range of scale and the distribution of persons at all scales is relatively uniform.

Density. In terms of density, on average there are ~ 28.87 persons per image in WiderPerson dataset, as shown in the fourth line of Table II. We also report the density from the existing datasets in Table III. Obviously, WiderPerson dataset is of much higher crowdness compared with all previous datasets. Caltech-USA suffers from extremely low-density, for that on average there is only ~ 1 person per image. The number in CityPersons reaches ~ 7 , a significant boost while still not dense enough. Both of them are insufficient to serve as an ideal benchmark for the challenging crowd scenes. Thanks to the pre-filtering and annotation protocol of our dataset, WiderPerson can reach a much better density. As shown in Table. II, we notice the density of persons are consistent across training/validation/testing subsets.

Diversity. Diversity is an important factor of a dataset. We compare the diversity of Caltech-USA, CityPersons and WiderPerson in Table. I. Since CityPersons testing set annotations are not publicly available, we only consider the training subset for a fair comparison. The Caltech-USA and KITTI datasets are recorded in one city at one season, and the CityPersons dataset is recorded across 18 cities, 3 countries

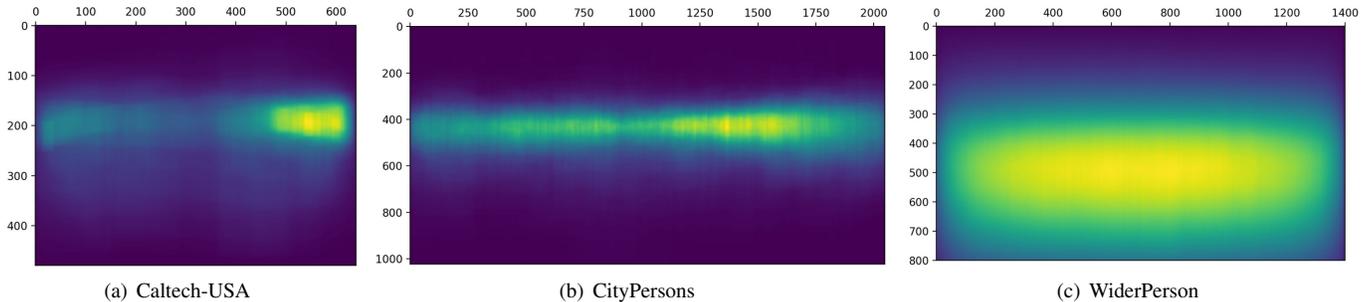


Fig. 6. The location distribution of pedestrians on the image. Pedestrians on the Caltech-USA and CityPersons dataset are distributed in a narrow band across the center of the image, while WiderPerson has an uniform location distribution.

and seasons, while our WiderPerson dataset has no limitations on these conditions. WiderPerson contains person in a wide range of scenarios, while Caltech-USA and CityPersons are all recorded by a car traversing on streets. In order to visualize the diversity of annotations on the different datasets, we count the location distribution of persons, *i.e.*, iterating over all person annotations, for each location, if it is inside one annotation, then its count plus 1. The images of Caltech and CityPersons have a fixed resolution (640×480 and 2048×1024 , respectively), while our dataset varies in size, so we resize all images into the same resolution (1400×800) to count the location distribution of persons. Fig. 6 shows the location distribution of persons for different dataset in the way of heat map. We can see that persons on the Caltech-USA and CityPersons dataset are distributed in a narrow band across the center of the image, *i.e.*, persons are concentrated on two sides of the road and mostly appear at the right side, since their images are collected by a biased data collection method that the car drives under the right-handed traffic condition. In contrast, our WiderPerson dataset has a uniform location distribution and persons appear in any position except the upper part (*i.e.*, the sky).

Also, the number of identical persons is another important evidence of diversity. As reported in the fifth line in Table I, the number of identical persons amounts up to $\sim 236k$ in our WiderPerson dataset. In contrast, the Caltech-USA dataset only contains $\sim 1,300$ unique pedestrians, since images in Caltech-USA are not sparsely sampled, resulting in less amount of identical persons. While CityPersons frames are sampled very sparsely and each person is considered as unique. Like CityPersons, each person on our dataset can be considered as unique, but one more order of magnitude. Besides, WiderPerson also provides fine-grained labels for persons. As shown in Fig. 7, pedestrians are the majority (64.8%). Partially-visible persons account for 29.9% since our dataset is dense. Although riders only occupy 0.6%, the absolute numbers are still considerable, as we have a large pool of $\sim 236k$ persons.

Occlusion. Occlusion is another important factor for evaluating the pedestrian detection performance. There are two types of occlusion: inter-class occlusion, which occurs when a person is occluded by stuff or objects from other categories; and intra-class occlusion (also referred to as crowd occlusion), which occurs when a person is occluded by other persons. As

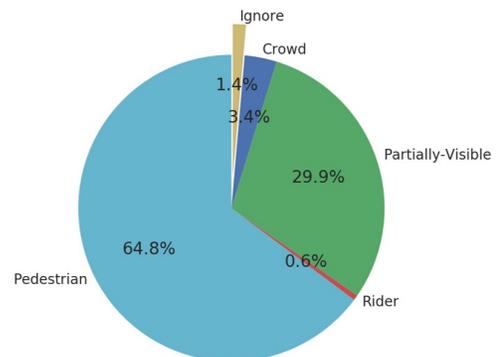


Fig. 7. Fine-grained person categories on WiderPerson.

TABLE IV
COMPARISON OF PAIR-WISE OVERLAP BETWEEN TWO PERSON INSTANCES.

pair/image	Caltech-USA	CityPersons	WiderPerson
IoU>0.3	0.06	0.96	9.21
IoU>0.4	0.03	0.58	4.78
IoU>0.5	0.02	0.32	2.15
IoU>0.6	0.01	0.17	0.81
IoU>0.7	0.00	0.08	0.24
IoU>0.8	0.00	0.02	0.06
IoU>0.9	0.00	0.00	0.01

described in [5], the intra-class occlusion is a more challenging issue than the inter-class occlusion. Lots of works have focused on the former problem and made great progress. However, the crowd occlusion has not been well researched and solved. On the one hand, it is difficult because of the problem itself. On the other hand, there is no suitable dataset. Therefore, we introduce this diversity and dense pedestrian detection dataset. To demonstrate its degree of crowd occlusion, we provide statistical information on pair-wise occlusion. For each image, we count the number of person pairs with different intersection over union (IoU) threshold. The results are shown in Table IV. In average, few person pairs with an IoU threshold of 0.3 are included in Caltech-USA. For CityPersons dataset, the number is less than one pair per image. However, the number is 9.21 for WiderPerson. Moreover, there are averagely 2.15 pairs whose IoU is greater than 0.5 in the WiderPerson dataset. These data can demonstrate that various occlusion levels are

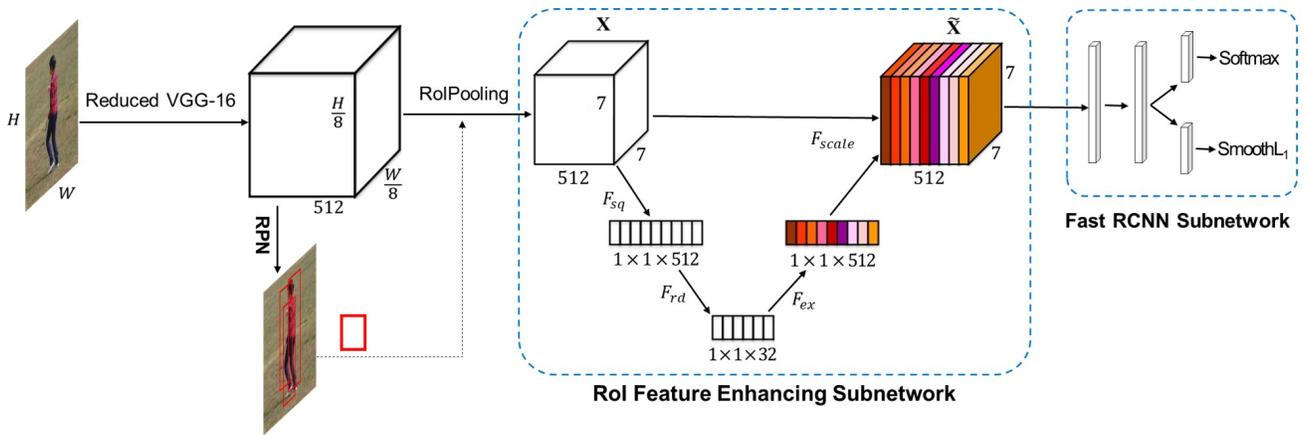


Fig. 8. Diagram of our improved Faster R-CNN. Reduced VGG-16 means removing the fourth max pooling and using the “hole algorithm”. RoI feature enhancing subnetwork is a reimplementation of SENet with an identical block structure, which consists of one global average pooling layer (F_{sq}) and two consecutive fully connected layers (F_{rd} and F_{ex}). RoI feature X is re-weighted to generate the output \tilde{X} of the SE block which then be fed directly into subsequent Fast R-CNN subnetwork.

well-represented in WiderPerson, especially heavily occluded cases, while they can be hardly found in previous datasets.

E. Benchmarking

With the publication of this paper, we will create a website for WiderPerson dataset, where its annotations for the training and validation subsets are made freely available to academic and non-profit organizations for non-commercial, scientific use. There is also an evaluation instruction on the website for researchers to evaluate the performance of their detectors over the held-out testing annotations. A leaderboard will be maintained and results are tallied online, either by name or anonymous.

We follow the same evaluation metric as used for Caltech-USA [1] and CityPersons [3], denoted as MR , which stands for the average log miss rate over false positives per-image ranging in $[10^{-2}, 10^0]$. MR is a suitable indicator for the algorithms applied in the real world applications. When evaluating pedestrian detection performance, riders/partially-visible persons/crowd/ignore regions are ignored, which means that those annotations are not considered as false negatives and detections matching with those annotations are not counted as false positives.

IV. PROVIDED BASELINE METHOD

Before delving into our new dataset, we first build two strong baseline detectors as a tool for our experiment analyses based on Faster R-CNN [11] and RetinaNet [58], which are two representative detectors from the two-stage and one-stage approach, respectively. We aim to find a straightforward architecture to provide good performance on WiderPerson.

A. Improved Faster R-CNN

Faster R-CNN is masterpieces of the detection framework for general object detection and has dominated the field of object detection in recent years. It essentially consists of two components: a fully convolutional Region Proposal Network

(RPN) for proposing candidate regions which likely contain objects, followed by a downstream Fast R-CNN network to classify a region of image into objects (and background) and refine the boundaries of those regions. Although competitive performance has been achieved on general object detection task, it under-performs on the pedestrian detection task (as reported in [42]). The reason behind its poor performance on pedestrian detection is that it fails to handle heavily occluded and dense pedestrians, which are dominant on our new dataset. In this work, we propose some improvements to extend the Faster RCNN architecture for occluded and dense pedestrian detection. Our improved Faster R-CNN is based on VGG-16 [59] and ResNet-50 [60] as they are very common in pedestrian detection [3], [42], [44]. We use the same anchor setting from [3], *i.e.*, 11 different anchor-box scales and 1 aspect ratio ($w/h = 0.41$) are used to capture objects across all sizes. Moreover, some improvements are proposed to boost the performance on pedestrian detection as follows.

Finer Feature Map. The vanilla Faster R-CNN uses a coarse feature map as the detection layer, *i.e.*, the last conv layer in the fifth block with stride of 16 pixels. Having such a coarse stride is harmful to small pedestrian detection, since it reduces the chances of having a high score over pedestrian and forces the network to handle large displacement relative to the object appearance. To increase the feature map resolution, we remove the fourth down-sampling operation and reduce the stride from 16 to 8 pixels, helping the detector to handle small pedestrians. Specifically, all layers before the fourth down-sampling operation are unchanged and all convolutional filters after it are modified by the “hole algorithm” [61] (*i.e.*, “Algorithm à trous”) to compensate for the reduced stride.

Ignore Region and Tiny Pedestrian Handling. We implement an ignore region handling for Faster R-CNN. Ignore regions might contain objects of a given class without precise localization. Simply treating these regions as background introduces confusing samples and has a negative impact on the detector quality. The ignore region handling prevents the

sampling of background boxes in those areas that could potentially overlap with real objects. Besides, training with very tiny samples could lead to models detecting a lot more false positives. Hence, we online filter pedestrians whose height is less than 20 pixels after scaling during training. Filtered pedestrians are handled as ignore regions in order to ensure that they are not sampled as background during training.

RoI Feature Enhancing. The RoIPooling layer uses max pooling to convert the features inside any valid region of interest into a fixed-size feature map, which is used by subsequent Fast R-CNN network to further classify and regress the proposals for final detections. Therefore, the representational ability of the pooled feature is the key to achieve high performance, especially on our highly diverse dataset. Inspired by [4], [62], we use a ‘‘Squeeze-and-Excitation’’ (SE) block to enhance the representational ability of the RoIPooling feature by explicitly modelling the interdependencies between the convolutional channels. More specific, the SE block performs sample-dependent feature re-weighting so as to select the more informative channel features while suppress less useful ones. As shown in Fig. 8, the newly added SE block is composed of one global average pooling layer and two consecutive fully connected layers, which is easy to implement and can obtain remarkable improvements while add little additional computational costs.

Dynamic Sample Strategy. The vanilla Faster R-CNN has a fixed sample strategy, *i.e.*, 256 and 128 samples for RPN and Fast R-CNN with 1 : 1 and 1 : 3 positive-negative ratio, respectively. Since there are ~ 28.87 persons per image in our dataset, the fixed sample strategy will lead to inadequate use of training positive samples. To solve this issue, we introduce a dynamic sample strategy: if there are too many positive samples, we determine the number of negative samples based on the above positive-negative ratio to ensure that all positive samples are used, otherwise we follow the original strategy.

B. Vanilla RetinaNet

In addition to the two-stage baseline detector, we also provide another baseline detector based on RetinaNet, the one-stage approach, which detects objects by regular and dense sampling over locations, scales and aspect ratios with high efficiency. RetinaNet proposes a focal loss to address the extreme foreground-background class imbalance by reshaping the standard cross entropy loss such that it down-weights the loss assigned to well-classified examples. We use the same setting of anchor scales as [58] and only modify the height *vs.* width ratio of anchors as 1:0.41 in consideration of the pedestrian shape.

V. EXPERIMENTS

In this section, we will introduce our implementation details about data processing and training setting. Notably, all the experiments are conducted based on the improved Faster R-CNN with VGG-16 unless otherwise specified. Firstly, we verify the effectiveness of our improvements via model analysis. Then, we conduct some experiments to analyse our WiderPerson

dataset in different aspects, including the detection result, quantity, quality and error. Finally, the generalization ability of our WiderPerson dataset will be evaluated on standard pedestrian benchmarks like Caltech-USA and CityPersons.

A. Implementation Detail.

Data Processing. To improve performance for small sized pedestrians, the input images are upscaled to a larger size using the bilinear interpolation algorithm. Specifically, the input image sizes of Caltech and CityPersons are set to $2\times$ and $1.3\times$ of the original images. As the images of WiderPerson are both collected from the Internet with various sizes, we resize the input so that their short edge is at 800 pixels while the long edge should be no more than 1400 pixels at the same time. We use horizontal image flipping as the only form of data augmentation. Multi-scale training and testing are not applied to ensure fair comparisons.

Training Setting. For the improved Faster R-CNN, all models are trained for $180k$ iterations with an initial learning rate of 0.01, and decreased by a factor of 10 after $120k$ on our WiderPerson dataset. On the CityPersons dataset, we set the learning rate to 10^{-3} for the first $40k$ iterations and decay it to 10^{-4} for another $20k$ iterations. On the Caltech-USA dataset, we train the network for $120k$ iterations with the initial learning rate 10^{-3} and decrease it by a factor of 10 after the first $80k$ iterations. To fine-tune the improved Faster R-CNN from WiderPerson to Caltech-USA and CityPersons, the number of iterations is the same but the learning rate is halved overall. All these models are optimized by the Stochastic Gradient Descent (SGD) algorithm on 1 TITAN X (Maxwell) GPU with a mini-batch 2. Weight decay and momentum are set to 0.0005 and 0.9. Besides, the RetinaNet baseline on the WiderPerson dataset is trained with 16 batch size for $25k$ iterations with 0.02 initial learning rate, which is then divided by 10 at $16k$ and again at $21k$ iterations.

B. Model Analysis

We carry out some ablation experiments on the WiderPerson validation subset to analyze our improved Faster R-CNN. For all the experiments, we use the same settings, except for specified changes to the components. We ablate each improvement one after another to examine how each proposed improvement

TABLE V
ANALYSIS OF PROPOSED IMPROVEMENTS. ALL MODELS ARE BASED ON FASTER R-CNN WITH VGG-16, TRAINED ON WIDERPERSON training SET AND TESTED ON validation SET. NUMBERS INDICATE *MR*.

Component	Step by step improvements					
	✓	✓	✓	✓	✓	
new anchor setting	✓	✓	✓	✓	✓	
finer feature map		✓	✓	✓	✓	
RoI feature enhancing			✓	✓	✓	
ignore region handling				✓	✓	
dynamic sample strategy					✓	
<i>Easy</i> subset	43.01	39.62	35.12	31.45	29.98	29.61
<i>Medium</i> subset	48.97	46.28	42.67	39.71	38.68	38.40
<i>Hard</i> subset	55.62	53.26	50.17	47.51	46.65	46.46

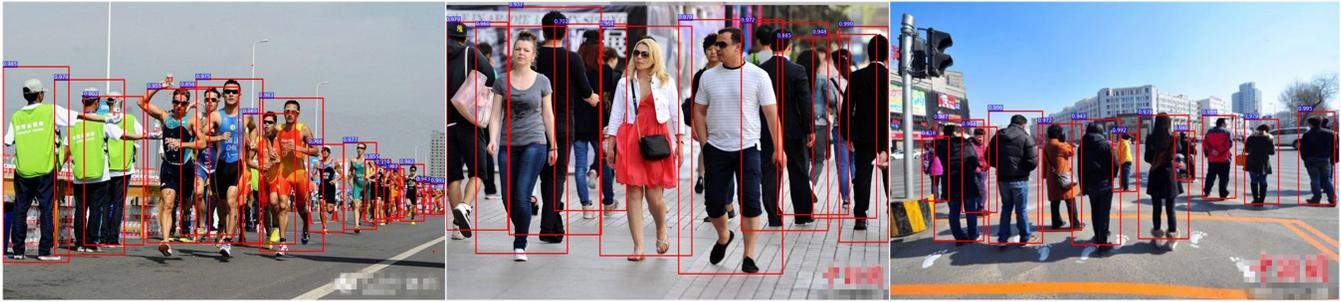


Fig. 9. Qualitative results for pedestrian detection of our improved Faster R-CNN with VGG-16 based on the WiderPerson dataset.

affects the final performance. Firstly, we replace the dynamic sample strategy with the original strategy. Secondly, the RoI feature enhancing module is ablated. Thirdly, we do not handle ignore regions and tiny ground truths during training stage. Fourthly, we do not reduce the VGG-16 backbone. Finally, we use original anchor scales rather than new anchor setting [3].

Some promising conclusions can be summed up according to the ablative results in Tab. V. Firstly, the new anchor setting is more suitable for the proposed dataset, which reduces the MR by 3.39%, 2.69% and 2.36% for Easy, Medium and Hard subset, respectively. Secondly, the finer feature map is used to provide more anchors and detailed information, which reduces the MR 39.62%, 46.28% and 53.26% to 35.12%, 42.67% and 50.17% for Easy, Medium and Hard subset, respectively, demonstrating its effectiveness. Thirdly, the ignore region and tiny pedestrian handling is proposed to ignore small ground truths and prevent the sampling of background boxes in those ignored areas. The comparison between the second and third columns in Tab. V demonstrates that it can bring 3.67% (Easy), 2.96% (Medium) and 2.66% (Hard) drops in MR, attributing to not involving confusing samples in training. Fourthly, according to the third and fourth columns, we can observe a drop in MR of 1.47% (Easy), 1.03% (Medium) and 0.86% (Hard), these sharp declines demonstrate the effectiveness of the RoI feature enhancing. Finally, the comparison between the fourth and fifth columns in Tab. V indicates that the dynamic sample strategy decreases the MR by 0.37% (Easy), 0.28% (Medium) and 0.19% (Hard), owing to making full use of training samples.

C. Dataset Analysis

All experiments in this subsection are trained based on WiderPerson training subset and the results are evaluated on the validation subset. Firstly, we evaluate our improved Faster R-CNN in detail on validation subset, then study on the quantity and quality, finally analyze common failure cases.

Detection Results. Table VI illustrates our baselines' results on the WiderPerson validation subset. On the one hand, we achieve promising MR performances, *i.e.*, 31.47%, 40.45%, 48.32% for the vanilla RetinaNet, 29.61%, 38.40%, 46.46% for the improved Faster R-CNN with VGG-16, and 28.75%, 37.82%, 46.06% for the improved Faster R-CNN with ResNet-50 on Easy, Medium and Hard subsets, respectively. On the other hand, from these results, we can find that the

TABLE VI
MR AND SPEED PERFORMANCE OF OUR BASELINES ON THE WIDERPERSON validation SUBSET.

Baseline	Backbone	FPS	Easy	Medium	Hard
Vanilla RetinaNet	ResNet-50	8.93	31.47	40.45	48.32
Improved FRCNN	VGG-16	0.83	29.61	38.40	46.46
Improved FRCNN	ResNet-50	0.77	28.75	37.82	46.06

proposed WiderPerson dataset is a challenging benchmark even for the state-of-the-art pedestrian detection algorithms. In Table VII and Table VIII, we also report detection results of the improved Faster R-CNN with VGG-16 on Caltech, *i.e.*, 5.49% MR, and CityPersons, *i.e.*, 12.49% MR. It further demonstrates that our WiderPerson dataset is much challenging than the standard pedestrian detection benchmarks based on the detection performance. Since our WiderPerson dataset varies largely in scenario and occlusion, which bring many difficulties to pedestrian detection. The illustrative examples of pedestrian detection based on our improved Faster R-CNN with VGG-16 are shown in Fig. 9.

Quantity Analysis. As indicated in [63], there is a logarithmic relation between the amount of training data and the performance of deep learning methods. To understand the impact of having a larger amount of training data, we show how the performance grows as training data increases on our

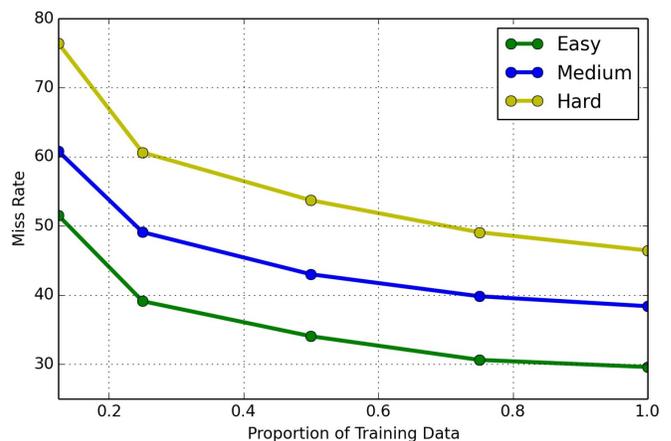


Fig. 10. Detection performance (MR) of our improved Faster R-CNN with VGG-16 as a function of training set size

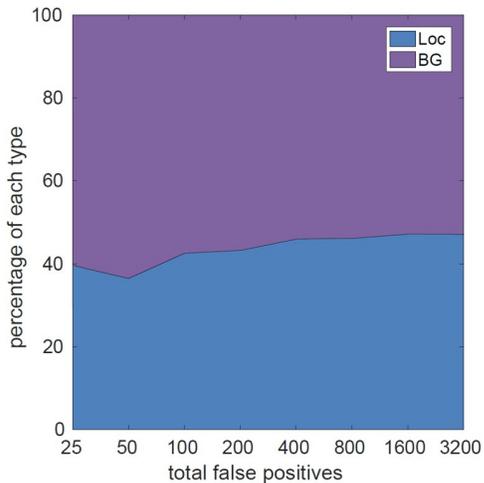


Fig. 11. Distribution of two error modes of false positives on the WiderPerson validation Hard subset.

benchmark. For this purpose, we train our baseline methods on different sized subsets which are randomly sampled from the training set. From Fig.10 we can observe that logarithmic relation between training set size and detection performance also holds on our benchmark for the improved Faster R-CNN across three subsets, *i.e.*, performance keeps improving with more data. Therefore, it is of great importance to provide CNNs with a large amount of data.

Quality Analysis. The importance of fine-grained annotations for riders and additional annotations for ignore regions is now examined. In ablation experiments, we have verified the effectiveness of ignored regions via training the model without ignore region handling, it is in accordance with earlier findings [3] that detection performance deteriorates when not using ignore regions during training. Besides, the evaluation protocol described in Section III-E ignores detected neighboring classes. For pedestrians this means that riders are not considered as false positives, hence training pedestrian detectors generally treat riders as ignore region. To verify whether rider annotations are useful for pedestrian detection, we can directly train the baseline detection method by including riders into pedestrians, since pedestrian and rider are annotated in the same way with a fixed aspect ratio on our dataset. As expected, comparing with training only with pedestrians, after these neighboring annotations are involved during training, detection performance increases on *MR* from 29.61%, 38.40% and 46.46% to 29.35%, 38.27% and 46.39% for Easy, Medium and Hard subset, respectively. Therefore, adding fine-grained annotations for riders is helpful for the pedestrian detection performance, since we can treat them as an additional training samples.

Error Analysis. We now utilize the detection analysis tool¹ to analyze the detection errors of our improved Faster R-CNN qualitatively on our WiderPerson validation dataset. The detection errors consist of false positives and false negatives. Firstly, we analyse the false positive errors. There are two

error modes of false positives in pedestrian detectors, *i.e.*, location (LOC) and background (BG). LOC indicates the localization errors that occurs when a pedestrian is detected with a misaligned bounding box, and BG indicates that a background region is mistakenly detected as a pedestrian. Fig. 11 shows the distribution of two types of false positives and BG seems the dominating error mode among top-scoring detection. Figure 12(a) illustrates some qualitative false positives of this method. As can be seen, animals, clothes and fakes are principal sources for confusion with real pedestrians. Certain pedestrian poses and aspect ratios can lead to multiple detections for the same pedestrian as shown in the *Multi Detections* category. Non-maximum suppression (NMS) is used by detection methods to suppress multiple detections. We use an *IoU* threshold of 0.5 which is not sufficient to suppress detections that have very diverse aspects.

Figure 12(b) illustrates some qualitative false negatives of this method. A lower *IoU* threshold would lead to more false negatives. These already occur for an *IoU* threshold of 0.5 as shown in the *NMS Repressing* category. Because of the high *IoU* between pedestrians, not all of them can be detected because of the greedy NMS. Thus, NMS is an important part of many deep learning methods that is usually not trained but has a great influence on detection performance. Small and occluded pedestrians are a further common source for false negatives. These two groups have also been analyzed in [64]. In some scenarios, usually only the lower part of a pedestrian is occluded due to various obstacles. In our qualitative analysis we have false negatives where the head is occluded. These are particularly challenging for pedestrian detection methods, as these cases are quite rare in the training dataset. Further challenges are unusual and extreme poses as shown in the *Others* group.

D. Generalization Capability

In this subsection, we evaluate the generalization capability of our WiderPerson dataset. As illustrated in Section III-D, the size of WiderPerson dataset is obviously more diverse and larger than the existing benchmarks, like Caltech-USA [1] and CityPersons [3]. Naturally, our dataset, with a reduced bias, should better capture the true world and result in superior generalization capabilities of the detectors which are trained on this dataset. To demonstrate the increased diversity of our dataset, we first train the model on our WiderPerson dataset and then fine-tune it on the other pedestrian detection benchmarks.

Caltech. The Caltech-USA dataset is one of the most popular and challenging datasets for pedestrian detection, which comes from approximately 10 hours 30Hz VGA video recorded by a car traversing the streets in the greater Los Angeles metropolitan area. We use the new high quality annotations provided by [65] to train and evaluate. The training and testing sets contains 42,782 and 4,024 frames, respectively. The results are shown for the Caltech-USA dataset in Table VII. The overall detection performance is superior for the cases in which WiderPerson is used for pre-training. Our improved Faster R-CNN achieves 5.49% *MR* for pedestrians on the

¹<http://web.engr.illinois.edu/~dhoiem/projects/detectionAnalysis>



(a) False positive: a background region is mistakenly detected as a pedestrian, or a pedestrian is detected with a misaligned bounding box.



(b) False negative: a pedestrian fails to be detected.

Fig. 12. Qualitative detection errors for our improved Faster R-CNN (green: true positives, red: false positives or false negatives).

Caltech-USA testing set for the reasonable setting. When we directly evaluate the Caltech-USA trained model on the proposed WiderPerson Easy subset, we get a very high MR of 82.79% since Caltech-USA has limited density and diversity. In contrast, our model trained on WiderPerson without fine-tuning achieves 9.72% MR and can be boost to 4.27% MR with fine-tuning. Based on the pre-training of our WiderPerson dataset, our algorithm has superior performance on the Caltech-USA benchmark against the one without WiderPerson pre-training, and performs on-par with the state-of-the-arts.

CityPersons. The CityPersons dataset is built upon the semantic segmentation dataset Cityscapes to provide a new dataset of interest for pedestrian detection. It is recorded across 18 different cities in Germany with 3 different seasons and various weather conditions. The dataset includes 5,000 images (2,975 for training, 500 for validation, and 1,525 for testing) with 35,000 manually annotated persons plus

TABLE VII
EXPERIMENTAL RESULTS ON CALTECH-USA.

Training	Testing	MR
Caltech-USA	Caltech-USA	5.49
Caltech-USA	WiderPerson (Easy)	82.79
WiderPerson	Caltech-USA	9.72
WiderPerson⇒Caltech-USA	Caltech-USA	4.27
	RPN+BF [42]	7.3
	HyperLearner [44]	5.5
	OR-RCNN [8]	4.1
	Repulsion Loss [5]	4.0

~ 13,000 ignore region annotations. Both the bounding boxes and visible parts of pedestrians are provided and there are approximately 7 pedestrians in average per image. The results are shown for the CityPersons dataset in Table VIII. The

TABLE VIII
EXPERIMENTAL RESULTS ON CITYPERSONS.

Training	Testing	<i>MR</i>
CityPersons	CityPersons	12.49
CityPersons	WiderPerson (Easy)	73.45
WiderPerson	CityPersons	16.17
WiderPerson⇒CityPersons	CityPersons	11.13
Adapted Faster RCNN [3]		12.8
Repulsion Loss [5]		11.6
OR-CNN [8]		11.0

same findings hold for the CityPersons benchmark. Training on CityPersons dataset and testing on WiderPerson Easy subset has 73.45% *MR*, while training on WiderPerson dataset and testing on CityPersons validation subset achieves 16.17% *MR*, indicating the difficulty and expandability of our dataset. Again, our improved Faster R-CNN model pre-trained on WiderPerson can reduce the *MR* from 12.49% to 11.13% that is on-par with the state-of-the-arts, demonstrating our WiderPerson dataset can serve as an effective pre-training dataset for pedestrian detection task.

Summary. The superior performance on both Caltech-USA and CityPersons datasets when using WiderPerson for pre-training indicates a high dataset diversity. Models trained on this dataset will have increased generalization capabilities. However, due to the dataset biases, solely training on a dataset from the other domain without fine-tuning results in worse detection performance. Despite the dataset biases, the models are able to learn general features for the task of pedestrians detection when pre-trained on WiderPerson which proves useful for other datasets as well after fine-tuning. Using transfer learning to pre-train a network on generic data and fine-tune on the target domain is widely applied and used to increase overall performance.

VI. CONCLUSION

Current pedestrian detection benchmark datasets have contributed to spurring interest and progress in pedestrian detection research. With the help of CNN, modern methods have achieved remarkable performance on these benchmarks. However, it is still difficult to assess for real world performance, since there is a gap in the diversity and density between existing pedestrian detection benchmarks and real world requirements: 1) most of current datasets are collected in the fixed traffic scenario, which significantly reduces the diversity of the foreground and background. 2) crowd scenarios with occluded pedestrian are still under represented, limiting the variations in density. These limitations have partially contributed to the failure of some algorithms in coping with heavy occlusion and atypical scenario. To move forward the field of pedestrian detection, we introduce a diverse and dense pedestrian detection dataset called WiderPerson, which consists of 13,382 images with 399,786 annotations and varies largely in scenario and occlusion. Providing high quality annotations, it enables new experiments both for training better models and as new test benchmark. We propose some strong

baseline detectors based on Faster R-CNN and RetinaNet to benchmark the state-of-the-art detector. The cross-dataset generalization results of WiderPerson dataset demonstrate that it is an effective training source for pedestrian detection and can help to achieve state-of-the-art performance on the Caltech-USA and CityPersons datasets. In the future, we will provide continuous improvements and additions to the WiderPerson dataset. Besides, we plan to annotate the head bounding box for each pedestrian and explore their relationship to facilitate further studies on the dense pedestrian detection.

ACKNOWLEDGMENT

This work was supported by the National Key Research and Development Plan (Grant No.2016YFC0801002), the Chinese National Natural Science Foundation Projects #61876179, #61872367, #61806203, Science and Technology Development Fund of Macau (No. 152/2017/A, 0025/2018/A1, 008/2019/A1). We also acknowledge the support of NVIDIA with the GPU donation for this research.

REFERENCES

- [1] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: An evaluation of the state of the art," *TPAMI*, pp. 743–761, 2012.
- [2] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the KITTI vision benchmark suite," in *CVPR*, 2012, pp. 3354–3361.
- [3] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," in *CVPR*, 2017, pp. 4457–4465.
- [4] S. Zhang, J. Yang, and B. Schiele, "Occluded pedestrian detection through guided attention in cnns," in *CVPR*, 2018.
- [5] X. Wang, T. Xiao, Y. Jiang, S. Shao, J. Sun, and C. Shen, "Repulsion loss: Detecting pedestrians in a crowd," in *CVPR*, 2018.
- [6] S. Wang, J. Cheng, H. Liu, F. Wang, and H. Zhou, "Pedestrian detection via body-part semantic and contextual information with dnn," *TMM*, 2018.
- [7] J. Li, X. Liang, S. Shen, T. Xu, J. Feng, and S. Yan, "Scale-aware fast R-CNN for pedestrian detection," *TMM*, vol. 20, no. 4, pp. 985–996, 2018.
- [8] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Occlusion-aware r-cnn: Detecting pedestrians in a crowd," in *ECCV*, 2018.
- [9] H. Idrees, I. Saleemi, C. Seibert, and M. Shah, "Multi-source multi-scale counting in extremely dense crowd images," in *CVPR*, 2013, pp. 2547–2554.
- [10] J. Shao, K. Kang, C. C. Loy, and X. Wang, "Deeply learned attributes for crowded scene understanding," in *CVPR*, 2015, pp. 4657–4666.
- [11] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: towards real-time object detection with region proposal networks," *TPAMI*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [12] S. Silberstein, D. Levi, V. Kogan, and R. Gazit, "Vision-based pedestrian detection for rear-view cameras," in *Intelligent Vehicles Symposium*, 2014, pp. 853–860.
- [13] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *CVPR*, 2005, pp. 886–893.
- [14] B. Wu and R. Nevatia, "Cluster boosted tree classifier for multi-view, multi-pose object detection," in *ICCV*, 2007, pp. 1–8.
- [15] A. Ess, B. Leibe, and L. J. V. Gool, "Depth and appearance for mobile scene analysis," in *ICCV*, 2007, pp. 1–8.
- [16] D. Gerónimo, A. Sappa, A. López, and D. Ponsa, "Adaptive image sampling and windows classification for on-board pedestrian detection," in *ICVS*, vol. 39, 2007.
- [17] G. Overett, L. Petersson, N. Brewer, L. Andersson, and N. Petersson, "A new pedestrian dataset for supervised learning," in *Intelligent Vehicles Symposium*, 2008, pp. 373–378.
- [18] M. Enzweiler and D. M. Gavrila, "Monocular pedestrian detection: Survey and experiments," *TPAMI*, vol. 31, no. 12, pp. 2179–2195, 2009.
- [19] C. Wojek, S. Walk, and B. Schiele, "Multi-cue onboard pedestrian detection," in *Computer Society Workshop on CVPR*, 2009, pp. 794–801.

- [20] X. Li, F. Flohr, Y. Yang, H. Xiong, M. Braun, S. Pan, K. Li, and D. M. Gavrila, "A new benchmark for vision-based cyclist detection," in *Intelligent Vehicles Symposium (IV)*, 2016, pp. 1028–1033.
- [21] S. Hwang, J. Park, N. Kim, Y. Choi, and I. S. Kweon, "Multispectral pedestrian detection: Benchmark dataset and baseline," in *CVPR*, 2015, pp. 1037–1045.
- [22] S. Zhang, R. Benenson, M. Omran, J. H. Hosang, and B. Schiele, "Towards reaching human performance in pedestrian detection," *TPAMI*, vol. 40, no. 4, pp. 973–986, 2018.
- [23] M. Braun, S. Krebs, F. Flohr, and D. M. Gavrila, "The eurocity persons dataset: A novel benchmark for object detection," *CoRR*, 2018.
- [24] A. Sundaresan and R. Chellappa, "Multicamera tracking of articulated human motion using shape and motion cues," *TIP*, vol. 18, no. 9, pp. 2114–2126, 2009.
- [25] X. Lan, S. Zhang, P. C. Yuen, and R. Chellappa, "Learning common and feature-specific patterns: A novel multiple-sparse-representation-based tracker," *TIP*, vol. 27, no. 4, pp. 2022–2037, 2018.
- [26] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg, "SSD: single shot multibox detector," in *ECCV*, 2016, pp. 21–37.
- [27] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Single-shot refinement neural network for object detection," in *CVPR*, 2018.
- [28] J. Li, Y. Wei, X. Liang, J. Dong, T. Xu, J. Feng, and S. Yan, "Attentive contexts for object detection," *TMM*, vol. 19, no. 5, pp. 944–954, 2017.
- [29] P. Dollár, Z. Tu, P. Perona, and S. J. Belongie, "Integral channel features," in *BMVC*, 2009, pp. 1–11.
- [30] J. Yan, X. Zhang, Z. Lei, S. Liao, and S. Z. Li, "Robust multi-resolution pedestrian detection in traffic scenes," in *CVPR*, 2013, pp. 3033–3040.
- [31] S. Zhang, C. Baukchage, and A. B. Cremers, "Informed haar-like features improve pedestrian detection," in *CVPR*, 2014, pp. 947–954.
- [32] P. Dollár, R. Appel, S. J. Belongie, and P. Perona, "Fast feature pyramids for object detection," *TPAMI*, vol. 36, no. 8, pp. 1532–1545, 2014.
- [33] S. Zhang, R. Benenson, and B. Schiele, "Filtered channel features for pedestrian detection," in *CVPR*, 2015, pp. 1751–1760.
- [34] S. Paisitkriangkrai, C. Shen, and A. van den Hengel, "Strengthening the effectiveness of pedestrian detection with spatially pooled features," in *ECCV*, 2014, pp. 546–561.
- [35] P. Sermanet, K. Kavukcuoglu, S. Chintala, and Y. LeCun, "Pedestrian detection with unsupervised multi-stage feature learning," in *CVPR*, 2013, pp. 3626–3633.
- [36] J. H. Hosang, M. Omran, R. Benenson, and B. Schiele, "Taking a deeper look at pedestrians," in *CVPR*, 2015, pp. 4073–4082.
- [37] Y. Tian, P. Luo, X. Wang, and X. Tang, "Pedestrian detection aided by deep learning semantic tasks," in *CVPR*, 2015, pp. 5079–5087.
- [38] G. Brazil, X. Yin, and X. Liu, "Illuminating pedestrians via simultaneous detection and segmentation," in *ICCV*, 2017, pp. 4960–4969.
- [39] Z. Cai, M. J. Saberian, and N. Vasconcelos, "Learning complexity-aware cascades for deep pedestrian detection," in *ICCV*, 2015, pp. 3361–3369.
- [40] A. Angelova, A. Krizhevsky, V. Vanhoucke, A. S. Ogale, and D. Ferguson, "Real-time pedestrian detection with deep network cascades," in *BMVC*, 2015, pp. 32.1–32.12.
- [41] F. Yang, W. Choi, and Y. Lin, "Exploit all the layers: Fast and accurate CNN object detector with scale dependent pooling and cascaded rejection classifiers," in *CVPR*, 2016.
- [42] L. Zhang, L. Lin, X. Liang, and K. He, "Is faster R-CNN doing well for pedestrian detection?" in *ECCV*, 2016, pp. 443–457.
- [43] Z. Cai, Q. Fan, R. S. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *ECCV*, 2016, pp. 354–370.
- [44] J. Mao, T. Xiao, Y. Jiang, and Z. Cao, "What can help pedestrian detection?" in *CVPR*, 2017, pp. 6034–6043.
- [45] W. Ouyang and X. Wang, "A discriminative deep model for pedestrian detection with occlusion handling," in *CVPR*, 2012, pp. 3258–3265.
- [46] Y. Tian, P. Luo, X. Wang, and X. Tang, "Deep learning strong parts for pedestrian detection," in *ICCV*, 2015, pp. 1904–1912.
- [47] C. Zhou and J. Yuan, "Learning to integrate occlusion-specific detectors for heavily occluded pedestrian detection," in *ACCV*, 2016, pp. 305–320.
- [48] B. Wu and R. Nevatia, "Detection of multiple, partially occluded humans in a single image by bayesian combination of edgelet part detectors," in *ICCV*, 2005, pp. 90–97.
- [49] G. Duan, H. Ai, and S. Lao, "A structural filter approach to human detection," in *ECCV*, 2010, pp. 238–251.
- [50] B. Leibe, E. Seemann, and B. Schiele, "Pedestrian detection in crowded scenes," in *CVPR*, 2005, pp. 878–885.
- [51] X. Wang, T. X. Han, and S. Yan, "An HOG-LBP human detector with partial occlusion handling," in *ICCV*, 2009, pp. 32–39.
- [52] W. Ouyang and X. Wang, "Single-pedestrian detection aided by multi-pedestrian detection," in *CVPR*, 2013, pp. 3198–3205.
- [53] B. Pepik, M. Stark, P. V. Gehler, and B. Schiele, "Occlusion patterns for object class detection," in *CVPR*, 2013, pp. 3286–3293.
- [54] C. Zhou and J. Yuan, "Multi-label learning of part detectors for heavily occluded pedestrian detection," in *ICCV*, 2017, pp. 3506–3515.
- [55] M. K. Mihçak and R. Venkatesan, "New iterative geometric methods for robust perceptual image hashing," in *Security and Privacy in Digital Rights Management, ACM Workshop*, 2001, pp. 13–21.
- [56] S. Yang, P. Luo, C. C. Loy, and X. Tang, "WIDER FACE: A face detection benchmark," in *CVPR*, 2016, pp. 5525–5533.
- [57] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *ECCV*, 2014, pp. 391–405.
- [58] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017.
- [59] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014.
- [60] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [61] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *CVPR*, 2015, pp. 3431–3440.
- [62] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR*, 2018.
- [63] C. Sun, A. Shrivastava, S. Singh, and A. Gupta, "Revisiting unreasonable effectiveness of data in deep learning era," in *ICCV*, 2017, pp. 843–852.
- [64] R. N. Rajaram, E. Ohn-Bar, and M. M. Trivedi, "An exploration of why and when pedestrian detection fails," in *ITSC*, 2015, pp. 2335–2340.
- [65] S. Zhang, R. Benenson, M. Omran, J. H. Hosang, and B. Schiele, "How far are we from solving pedestrian detection?" in *CVPR*, 2016.