

Fused Deep Neural Networks for Efficient Pedestrian Detection

Xianzhi Du, Mostafa El-Khamy, Vlad I. Morariu, Jungwon Lee, and Larry Davis

xianzhi@umiacs.umd.edu



Abstract—In this paper, we present an efficient pedestrian detection system, designed by fusion of multiple deep neural network (DNN) systems. Pedestrian candidates are first generated by a single shot convolutional multi-box detector at different locations with various scales and aspect ratios. The candidate generator is designed to provide the majority of ground truth pedestrian annotations at the cost of a large number of false positives. Then, a classification system using the idea of ensemble learning is deployed to improve the detection accuracy. The classification system further classifies the generated candidates based on opinions of multiple deep verification networks and a fusion network which utilizes a novel soft-rejection fusion method to adjust the confidence in the detection results. To improve the training of the deep verification networks, a novel soft-label method is devised to assign floating point labels to the generated pedestrian candidates. A deep context aggregation semantic segmentation network also provides pixel-level classification of the scene and its results are softly fused with the detection results by the single shot detector. Our pedestrian detector compared favorably to state-of-art methods on all popular pedestrian detection datasets. For example, our fused DNN has better detection accuracy on the Caltech Pedestrian dataset than all previous state of art methods, while also being the fastest. We significantly improved the log-average miss rate on the Caltech pedestrian dataset to 7.67% and achieved the new state-of-the-art.

Index Terms—Pedestrian detection, object detection, semantic segmentation, soft-label, soft-rejection, network fusion, ensemble learning, convolutional neural network.

1 INTRODUCTION

Object detection is an essential problem in computer vision which aims to detect the locations of semantic objects in videos or digital images. Object detection is widely used in areas such as image retrieval, object identification, video surveillance, etc. Pedestrian detection is a branch of object detection problem which deals with detecting the specific human class. It has applications in various topics such as advanced driver assistance systems, person identification, face recognition, etc.

The pedestrian detection problem can be decomposed into region proposal generation, feature extraction, and pedestrian verification. In general, object detection involves generating candidates for bounding boxes enclosing the objects of interest, extracting robust features as high level representations of the candidates, and verifying each candidate to be a true or a false positive. In recent years, convolutional neural network based techniques have successfully been applied to pedestrian detection and achieved better performances in many challenging scenarios. Li et al. [1] trained multiple Fast R-CNN [2] based networks to detect

pedestrians with different scales and combined results from all networks to generate the final results. Hosang et al. [3] used the SquaresChnFtrs [4] method to generate pedestrian proposals and trained AlexNet [5] to perform pedestrian verification. Zhang et al. [6] used a Region Proposal Network (RPN) [7] to compute pedestrian candidates and a cascaded Boosted Forest [8] to perform sample re-weighting to classify the candidates.

In this paper, we propose a deep neural network fusion architecture to address the pedestrian detection problem, called Fused Deep Neural Network. Compared to previous methods, our proposed system is faster while achieving better detection accuracy. The architecture consists of a pedestrian candidate generator, which is obtained by training a deep convolutional neural network trained as a single shot detector (SSD) with a high detection rate, albeit a large false positive rate. A novel soft-rejection strategy is used to adjust the confidence in the detector candidates by fusion with a classification network employing an ensemble learning approach, and semantic segmentation network which provides pixel-wise labeling. The classification network deploys an ensemble of deep neural networks trained independently as verification networks, and their results are softly fused together with the detection results using the novel soft-rejection method. To prepare the training data for the verification networks, we devise a novel soft-label method to assign floating point labels to the detected candidates. Unlike traditional hard-label method for object verification, where binary labels are used, the value of the our soft-label is set to be the largest overlap ratio between the current detected bounding box and all the ground-truth bounding boxes, and is adjusted to saturate to the binary values. Additional accuracy improvements can be achieved at the expense of speed by the parallel semantic segmentation network which provides pixel-wise labels to generate a segmentation mask that delivers another soft confidence vote on the generated pedestrian candidates, and are further fused within the soft fusion framework. The proposed network architecture is shown in Figure 1.

Some of the ideas in this paper were presented at the 2017 IEEE WACV [9]. In this paper, we provide more detailed analysis of these ideas, show results on more datasets, and provide additional enhancements that improved the performance over that presented at the 2017 IEEE WACV [9]. The new techniques which we present here and helped provide the additional gains over [9] are the soft-label method for training classification methods, learning the parameters of soft-rejection fusion by an additional fusion network, and the new kernel based method to fuse the results of the semantic segmentation system and the detection

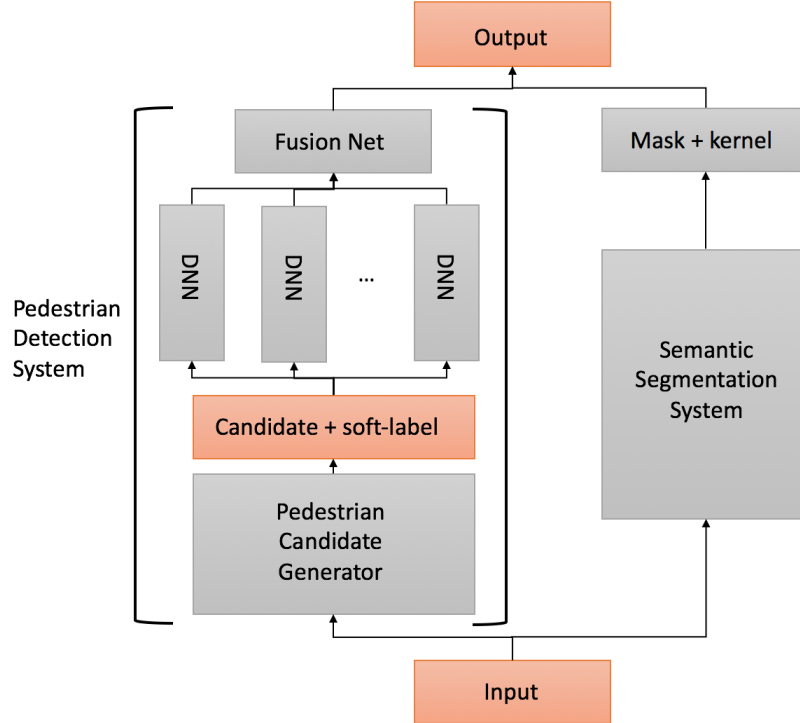


Fig. 1: Proposed Fused DNN Architecture

system. The new techniques of this paper helped to significantly increase the detection accuracy on the Caltech dataset from 8.18% to 7.65%. We also extend the model to work on more classes besides pedestrians, such as cars and cyclists. Besides the Caltech Pedestrian dataset, we evaluated on three more popular pedestrian detection datasets: INRIA, ETH, and KITTI. Our techniques performed better than all the previous state-of-the-arts on Caltech, INRIA, and ETH in both accuracy and speed, and achieved comparable results on KITTI. More ablation analysis is conducted to explain the effectiveness of our system.

The rest of this paper is organized as follows. Section 2 provides a detailed description of the pedestrian detection system. Section 3 describes the semantic segmentation system and how it helps to refine the detection results. Section 4 discusses the experiment results and explores the effectiveness of each component of the system. Section 5 draws conclusions on this work.

2 PEDESTRIAN DETECTION SYSTEM

2.1 Pedestrian Candidate Generator

In order to quickly obtain pedestrian candidates in various sizes and aspect ratios at all possible locations of the input image, we use an single shot multi-box detector (SSD) [10] as the candidate generator. The main reason we select SSD instead of other systems is that it uses multiple feature maps as the output layers. By lowering the accepting threshold of the confidence score, it outputs a large number of pedestrian candidates which are very likely to cover all the ground-truth pedestrians. Since it has a fully convolutional framework, it's also very fast.

The SSD network is a feed-forward convolutional network which consists of a VGG16 base, 8 convolutional layers above it, and a global pooling layer at the top. In the VGG16 base, the kernel size of the pool5 layer is set to 3×3 and the stride is

set to one, and the fc6 and fc7 layers are converted to dilated convolutional layers. Bounding box regression and classification are performed on the feature maps of 'conv4_3', 'fc7', 'conv6_2', 'conv7_2', 'conv8_2', 'conv9_2', and 'pool6' to generate the pedestrian candidates. Since the output of 'conv4_3' has a much larger feature scale than that of other filters, an L_2 normalization technique is used to scale down the feature magnitudes [11]. After each output layer, bounding box (BB) regression and classification are performed to generate pedestrian candidates. The network architecture of the deployed pedestrian SSD is shown in Figure 2.

To generate the pedestrian candidates, a set of default bounding boxes are slid on top of each output feature map. For each output layer of size $m \times n \times p$, a set of default BBs at different scales and aspect ratios are placed at each location. At every pixel location of the 7 output layers, we place 6 default bounding boxes (BBs) with aspect ratios [0.1, 0.2, 0.41, 0.8, 1.6, 3.0] and relative heights [0.05, 0.1, 0.24, 0.38, 0.52, 0.66, 0.80]. Since 0.41 is the average aspect ratio of all the pedestrian annotations, we place another set of default bounding boxes with that aspect ratio and relative heights [0.1, 0.24, 0.38, 0.52, 0.66, 0.80, 0.94]. In the training stage, we further categorize all the pedestrians into three classes: 'Full pedestrian', 'Occluded pedestrian', and 'People'. The 'People' class is defined as a group of people that are very close to or overlap with each other. The aspect ratios 0.6 and 3.0 are designed for 'People'. At each default bounding box, a $3 \times 3 \times p$ convolutional kernel is applied to produce classification scores, and perform bounding box regression by calculating the location offsets with respect to the default bounding box.

The multi-task training objective of SSD is given by Equation (1)

$$L = \frac{1}{N}(L_{conf} + \alpha L_{loc}) \quad (1)$$

where L_{conf} is the softmax loss for the classification task and

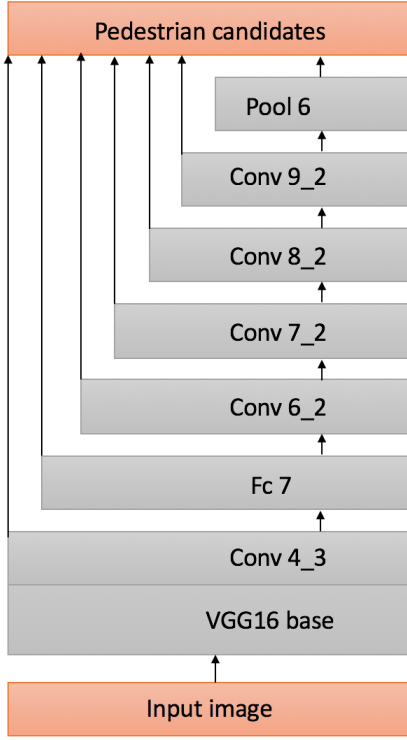


Fig. 2: Fully Convolution

L_{loc} is the smooth L_1 localization loss [2], N is the number of positive default boxes, and α is a constant weighting factor to keep a balance between the two losses. Since SSD uses 7 output layers to generate multi-scale BB outputs, it provides a large pool of pedestrian candidates varying in scales and aspect ratios.

When training the SSD as the primary candidate generator, a default detection BB is labeled as positive if it has a Jaccard overlap (intersection over union ratio) greater than 0.5 with any ground truth BB, otherwise it is labeled as negative. The SSD generated pedestrian candidates should cover almost all ground truth pedestrians, while being as accurate as possible. This is essential to our fused DNN architecture, where following classification networks will attempt to improve the precision by decreasing the confidence in many of the false positives introduced by the SSD, while preserving the recall rate. By being a fully convolutional network, SSD is a fast candidate generator.

2.2 Classification System

2.2.1 Classification networks with soft-label method

In this section, we describe the soft-label method we devised for preparing the training data for the classification network. The hard-label used in common object detection networks assigns a binary label to each pedestrian candidate bounding box by thresholding the Jaccard overlap ratio between this bounding box and the ground-truth bounding boxes. However, this is not the optimal strategy, especially when the overlap ratio is close to the threshold. In this work, we introduce the soft-label method to label pedestrian candidates for further classification. The soft-label method will assign a floating point label to each pedestrian candidate using the largest overlap ratio between the current pedestrian candidate and all the ground-truth bounding boxes. Suppose we have a pedestrian candidate and the ground-truth

bounding box it overlaps with the most. The soft-labels for the pedestrian class label $label_{ped}$ and the background class label $label_{bg}$ are respectively calculated by Equations (2) and (3), where A_{BB_d} and A_{BB_g} represent the areas covered by the detection BB and the ground truth BB, respectively

$$label_{ped} = \frac{A_{BB_d} \cap A_{BB_g}}{A_{BB_d}} \quad (2)$$

$$label_{bg} = 1 - \frac{A_{BB_d} \cap A_{BB_g}}{A_{BB_d}} \quad (3)$$

However, we found that it is better to use a soft-label in cases of mild confidence, and use a hard-label in cases when the confidence in the label exceeds a threshold. Hence, we devised a hybrid soft-label method as follows: If the overlap ratio between the pedestrian candidate and the ground truth bounding box is lower than a threshold th_a or greater than a threshold th_b , we think it's safe to label the current sample as background or pedestrian with probability one. For the candidates with intermediate confidence, the soft-label method is used to assign floating point labels, and we normalize the range $[th_a, th_b]$ to $[0, 1]$. Moreover, for convenience of implementation, the soft-label is made continuous over its range from 0 to 1. The final soft-label method is shown by Equations (4) and (5).

$$label_{ped} = \begin{cases} 1, & \text{if } \frac{A_{BB_d} \cap A_{BB_g}}{A_{BB_d}} > th_b \\ 0, & \text{if } \frac{A_{BB_d} \cap A_{BB_g}}{A_{BB_d}} < th_a \\ \frac{\frac{A_{BB_d} \cap A_{BB_g}}{A_{BB_d}} - th_a}{th_b - th_a}, & \text{otherwise,} \end{cases} \quad (4)$$

$$label_{bg} = 1 - label_{ped}. \quad (5)$$

The classification networks are trained to minimize the cross entropy loss objective function,

$$\epsilon = - \sum_{i=1}^c l_i \log y_i \quad (6)$$

where c is the number of classes, l_i and y_i are the soft-label and the softmax probability for class i . Note that $\sum_i l_i = 1$. Note that for the conventional binary labeling method, l_i is the indicator function, which is 1 for the correct class and 0 otherwise. Minimizing the cross entropy is equivalent to maximizing the softmax probability of the correct class. In our soft-label case, the softmax probabilities of all the classes are used to contribute to the cross entropy loss. The floating point soft-labels will determine how much each class contributes. When doing back-propagation, the derivative of the cross entropy cost function with respect to class i is calculated as Equation (7).

$$\begin{aligned} \frac{\partial \epsilon}{\partial z_i} &= - \sum_{j=1}^c l_j \frac{\partial \log y_j}{\partial z_i} = - \frac{l_i}{y_i} \frac{\partial y_i}{\partial z_i} - \sum_{j \neq i} \frac{l_j}{y_j} \frac{\partial y_j}{\partial z_i} \\ &= - \frac{l_i}{y_i} y_i (1 - y_i) - \sum_{j \neq i} \frac{l_j}{y_j} (-y_j y_i) \\ &= -l_i + y_i \sum_{j=1}^c l_j = y_i - l_i \end{aligned} \quad (7)$$

We note that Equation (7) shows the gradient calculations and holds for both conventional hard-label method as well as for the

proposed soft-label method. For conventional hard-label method, a training sample has label 1 for the correct class, and label 0 for the incorrect classes. We note that it holds for the soft-label method as long as the sum of the soft-labels over all classes is 1.

As shown in Figure 3, the classification network constitutes of multiple networks, run in parallel, where we used the idea of ensemble learning. The constituent networks are deep classification neural networks which have different network structures, but trained with the same input data. All pedestrian candidates generated by the primary candidate generator with confidence score greater than 0.01 and height greater than 40 pixels are collected as the new training data for the classification network. To solve the problem of classifying candidates with different aspect ratios and sizes, all candidate BBs are rescaled to a fixed input size. The goal of these secondary classification networks is to further classify the detections from the first stage single shot detector as a true detection or a false detection.

2.2.2 Soft-rejection based fusion

The opinions of all the constituent classification networks are fused with those of the candidate generator (CG). By doing so, it's more likely to get a lower error than having a single network. Since it's hard to bias towards each of the single network, it's also less likely to over-fit. Also, since the networks are run in parallel, the speed of the classification network is limited by its slowest constituent network. There are several conventional methods commonly used for opinion fusion, such as computing the mean of all results, majority voting, or the hard-rejection method. The hard-rejection method will eliminate a pedestrian candidate based on a single negative vote from one classification network. Instead, we introduce the soft-rejection network fusion (SNF) method, as we use classification networks with different structures, and we expect each network to work well in some of the subcategories while performing mediocre in other categories. The soft-rejection based fusion method can be described as follows: For one pedestrian candidate, the k_{th} classification network gives us a classification probability p_k . If p_k is higher than a threshold t_1 , we generate a scaling factor s_k greater than one to boost the initial confidence score generated by the SSD candidate generator. Otherwise, we generate a scaling factor less than one to decrease the initial confidence score. To prevent any of the classification networks from dominating the final results, we set a lower bound t_2 to the scaling factors. The scaling factors coming out from all the classification networks are further multiplied together with the initial confidence score to produce the final score. The idea behind this is that instead of accepting or rejecting any candidate, we boost or decrease their scores instead. This is because a poor classification network can be compensated by other good ones with SNF, whereas a wrong elimination of a true pedestrian by hard-rejection cannot be corrected. The SNF idea is illustrated in Equation (8) and Equation (9).

$$s_k = \max\left(p_k \times \frac{1}{t_1}, t_2\right) \quad (8)$$

$$S_{FDNN} = S_{CG} \times \prod_{k=1}^K s_k. \quad (9)$$

2.2.3 Soft-rejection based fusion network

The values of all the parameters in the soft-rejection based fusion method as described in Subsection 2.2.2 were selected by cross-validation in our previous work [9]. However, we found that the

parameters when optimized on one dataset do not easily generalize to other datasets. Instead of hand-crafting such parameters, we propose in this paper to use a neural network to learn the optimal parameters, while keeping the idea of the soft-rejection based fusion method. We call this new method soft-rejection based fusion network.

Let p_1, p_2, \dots, p_K be the inputs to the fusion network, where p_k is the softmax output of the k_{th} classification network. The input layer also deploys a log layer to get the classification log probabilities of the individual networks. The input layer is followed by two fully connected layers, each has 500 neurons, and one softmax layer to predict the weights for each classification network. The output of the neural network is the exponent of the weighted sum of all classification log probabilities. This results in a soft-fusion scheme which scales the candidate generator's confidence scores with the product of (exponentially) weighted classification probabilities of all individual classification networks, where the weights w_k are optimized by the fusion network, and adheres to our soft network fusion framework described by Equation (9),

$$\begin{aligned} S_{FDNN} &= S_{CG} * \exp\left(\sum_{k=1}^K w_k * \log(p_k)\right) \\ &= S_{CG} * \prod_{k=1}^K p_k^{w_k}. \end{aligned} \quad (10)$$

2.2.4 Training the classification system

There are two ways to train the classification system. The first method is to train an end-to-end system. For all the classification networks, we remove their loss layers and concatenate the output neurons for the pedestrian class from the softmax layers to form the input layer to the fusion network. This system has classification networks as branches and merged together by the fusion network at the end. This is shown at the left of Figure 3. However, since we train all the networks together, the structure grows huge and is prone to overfitting. What's more, since all the branches have different structures, they require different settings of optimal hyper-parameters and have different converging speed.

The second method is to train the classification networks first, and use the output probabilities as inputs to train the fusion network separately. Since this is straightforward and easy to implement, as shown at the right figure of Figure 3, this method is finally used in this paper.

3 PIXEL-WISE SEMANTIC SEGMENTATION SYSTEM

Based on the idea of fusion of multiple experts, we propose to enhance the accuracy of the system by adding an expert network which is trained independently from the candidate generator. In this implementation, we chose the independent network to be a semantic segmentation network trained to provide pixel-wise classification of pixels, in contrast to the region classification done by the candidate generator. In our work, we utilize a semantic segmentation (SS) network based on deep dilated convolutions and context aggregation [12] running in parallel with the pedestrian detection system to further refine the end detection results of the whole system. The SS network is trained on the Cityscapes dataset for driving scene segmentation [13]. To perform dense prediction, the network consists of a fully convolutional VGG16 network,

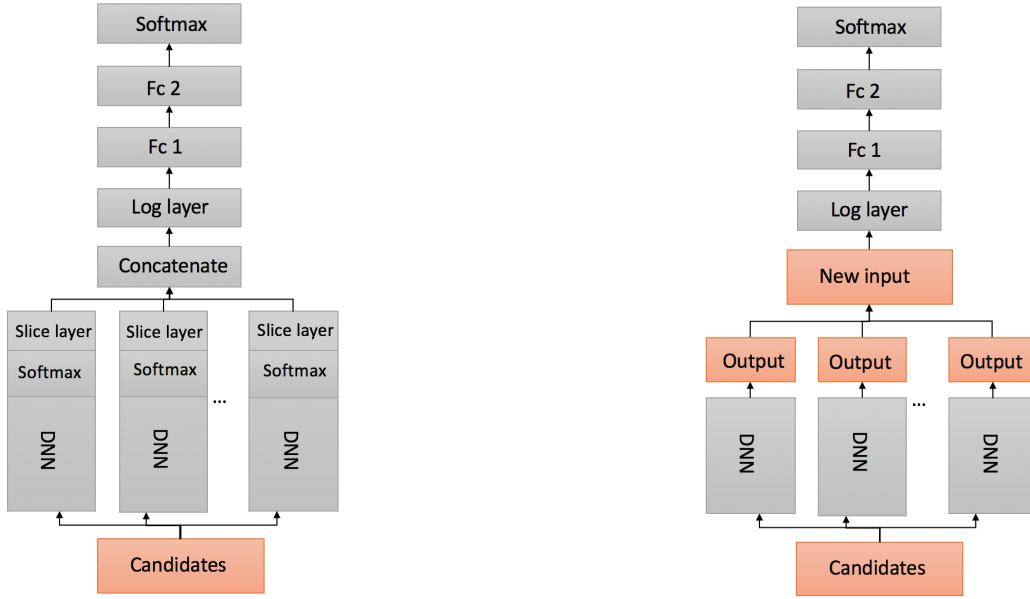


Fig. 3: The two fusion network designs described in the paper. The left structure is an end-to-end training scheme. The right structure trains all the networks separately.

adapted with dilated convolutions as the front end prediction module, whose output is fed to a multi-scale context aggregation module, consisting of a fully convolutional network whose convolutional layers have increasing dilation factors. We consider both the "person" and "rider" classes in Cityscapes dataset as pedestrians, and the remaining classes as background.

3.1 Soft Fusion with a Pixel-wise Semantic Segmentation System

The results of the semantic segmentation system are fused with those of the pedestrian detection system as follows: First, we process the same original input image using the semantic segmentation system. This produces a binary mask where the foreground pixels are set to 1 to represent the class of interest (e.g. pedestrian) and the background pixels are set to 0. Then, for each of the bounding boxes detected produced by the candidate generator, we analysis its pixels at the same locations on the binary mask. The soft fusion scaling factor is computed as the weighted sum of the foreground pixels within the bounding box, where the weighting factors are calculated from a weight matrix, called the Kernel, within the bounding box. To enable this weighted sum, all detected bounding boxes and their overlapped masks are rescaled to have the same size as the Kernel. The Kernel is trained as the mean of the semantic segmentation binary masks of all ground-truth pedestrian bounding boxes in the training set and normalized to have a sum of 1. This fusion method is illustrated in Equations (11) and (12).

$$S_{SS} = \sum_{i,j}^{A_{BB}} \text{mask}(i,j) \times \text{Kernel}(i,j), \quad (11)$$

$$S_{FDNN2+SS} = S_{FDNN2} \times S_{SS} \quad (12)$$

where the A_{BB} is the area of the bounding box. $\text{mask}(i,j)$ and $\text{Kernel}(i,j)$ are the pixel value of the binary mask and the Kernel at location (i,j) . From the visualization of the weight mask,

illustrated in Figure 4, we see that the pixels at the center of the kernel tend to have higher values than the pixels at the boundary. This agrees with the fact that in a perfect detection, the object of interest tends to appear at the center of the bounding box. We can see that the Kernel will have the effect of boosting the score of a detection whose bounding box fits the object of interest (e.g. pedestrian) and decreasing the score of a detection whose bounding box is not well located. This model will be referred to as 'FDNN2 + SS'

It is worth noting that in our previous work [9], the soft fusion of the semantic segmentation results with that of the object detection system, labeled 'FDNN + SS', was done differently. For the sake of completeness and comparison of the results, we describe it here. The SS mask is intersected with all detected BBs produced by the CG and the degree to which each candidate's BB overlaps with the pedestrian category in the SS activation mask gives a measure of the confidence of the SS network in the candidate generator's results. If the pedestrian pixels occupy at least 20% of the candidate BB area, we accept the candidate and keep its score unaltered; Otherwise, the candidate generator's scores are softly fused by scaling as in Equation 13,

$$S_{FDNN+SS} = \begin{cases} S_{FDNN}, & \text{if } \frac{A_M}{A_{BB}} > 0.2 \\ S_{FDNN} \times \max(\frac{A_M}{A_{BB}} \times a_{ss}, b_{ss}), & \text{otherwise} \end{cases} \quad (13)$$

where A_{BB} represents the area of the BB, A_M represents the area within A_{BB} covered by the semantic segmentation mask, a_{ss} , and b_{ss} are chosen as 4 and 0.35 by cross validation.

We also note that SNF of the CG network with an SS network is slightly different from SNF with a classification network. The reason is that the SS network can generate new detections which have not been produced by the CG, which is not the case for the classification networks. To address this, the proposed SNF methods 'FDNN2 + SS' and 'FDNN + SS' eliminate new detections from the SS network, which do not overlap with any CG detection.

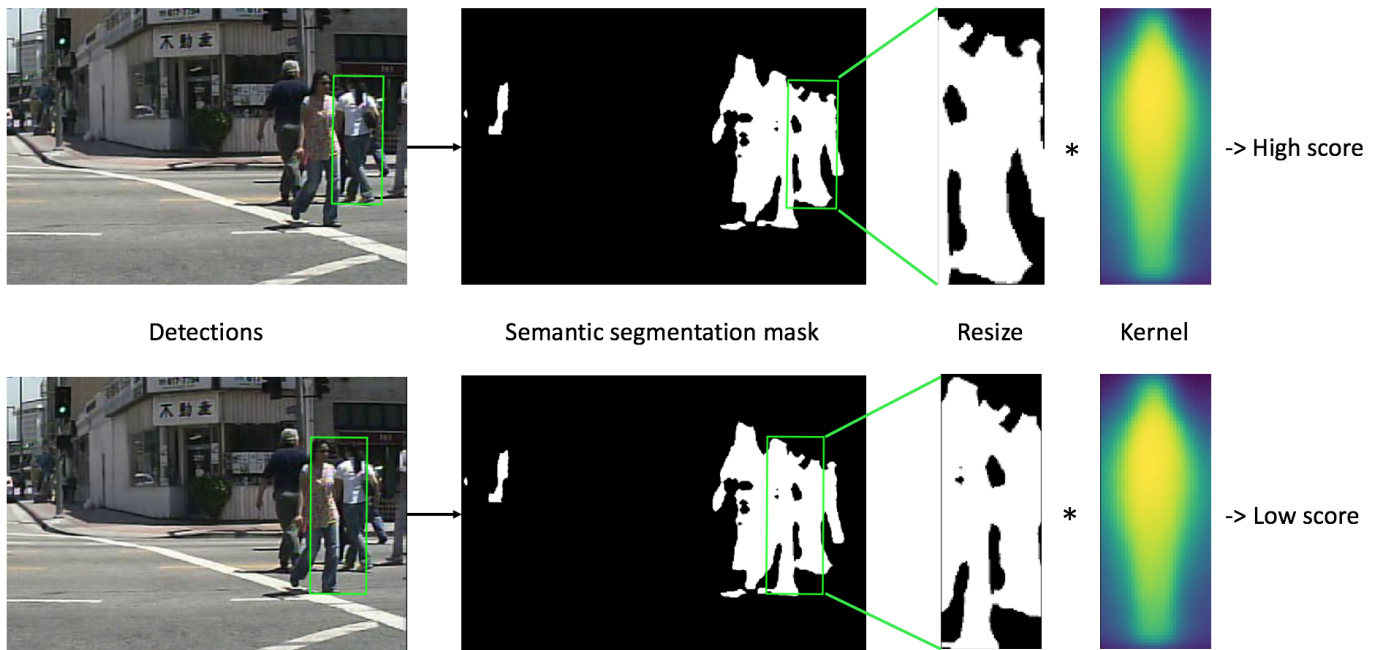


Fig. 4: The idea of kernel-based method to fuse the semantic segmentation system and the detection system.

4 EXPERIMENTS AND RESULT ANALYSIS

4.1 Training settings

The proposed method is trained on the training sets of the Caltech Pedestrian dataset, the ETH dataset, and the TudBrussels dataset.

To train the pedestrian candidate generator, both the original images and the horizontally flipped images which contain at least one annotated bounding box are used, which results in around 68,000 training images in total. Among all the annotated bounding boxes, there are about 109,000 annotated bounding boxes in 'Person_full' class, 60,000 annotated bounding boxes in 'Person_occluded' class, and 35,000 bounding boxes in 'People' class. All the images are of size 480×640 . The model is fine-tuned from the Microsoft COCO [14] pre-trained SSD model for 40,000 iterations using the standard stochastic gradient descent (SGD) algorithm and the back-propagation algorithm at a learning rate of 10^{-5} .

To train the classification system, all the ground-truth annotations and the pedestrian candidates generated from the previous stage with height greater than 40 pixels and confidence score larger than 0.01 are selected, and rescaled into a fixed size of 250×250 to represent the training samples. For data augmentation, a 224×224 patch is randomly cropped out of each training sample and horizontally flipped with probability 0.5. To label the training samples, the soft-label method as described by Equations (4) and (5) is implemented. The thresholds th_a and th_b are set to 0.4 and 0.6, respectively. To build the classification networks, one ResNet-50 [15] and one GoogleNet [16] are used as the classification networks. Both of the classifiers are fine-tuned from the ImageNet pre-trained models using the standard SGD algorithm and the back-propagation algorithm at a learning rate of 10^{-4} .

As we don't have pixel-level labels for pedestrian detection datasets to incorporate the semantic segmentation network, the dilated convolution model [12] trained on the Cityscapes dataset

[13] is directly implemented. All the classes are considered as background, except the 'Person' and 'Rider' classes which are considered as the 'Pedestrian' class. Due to the lack of well-labeled pedestrian datasets for our problem, no fine-tuning is involved in this step. All the input images are rescaled from 480×640 into 1024×2048 . To preserve the aspect ratio so as to preserve the human body shape, the image's height is firstly scaled to 1024 and then blank patches are padded on both left and right sides.

All the above mentioned models are built with the Caffe deep learning framework [17].

4.2 Evaluation settings and results

We evaluate the proposed method on the four most popular pedestrian detection datasets: the Caltech Pedestrian dataset, the INRIA dataset, the ETH dataset, and the KITTI dataset. The log-average miss rate (L-AMR) is used as the performance evaluation metric [18]. L-AMR is computed evenly spaced in log-space in the range 10^{-2} to 10^0 by averaging the miss rate at the rate of nine false positives per image (FPPI) [18]. There are multiple evaluation settings defined based on the height and visible part of the bounding boxes. The most popular settings are listed in Table 1.

Setting	Description
Reasonable	50+ pixels. Occ. none or partial
All	20+ pixels. Occ. none, partial, or heavy
Far	30- pixels
Medium	30-80 pixels
Near	80+ pixels
Occ. none	0% occluded
Occ. partial	1-35% occluded
Occ. heavy	35-80% occluded

TABLE 1: Evaluation settings for Caltech Pedestrian dataset.

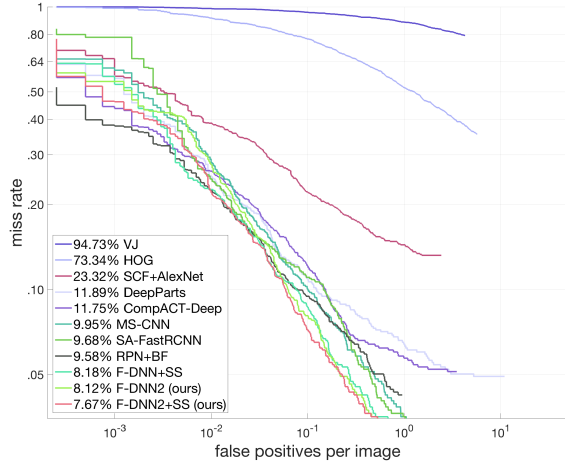


Fig. 5: L-AMR vs. FPPI plot under the 'Reasonable' evaluation setting on Caltech Pedestrian dataset.

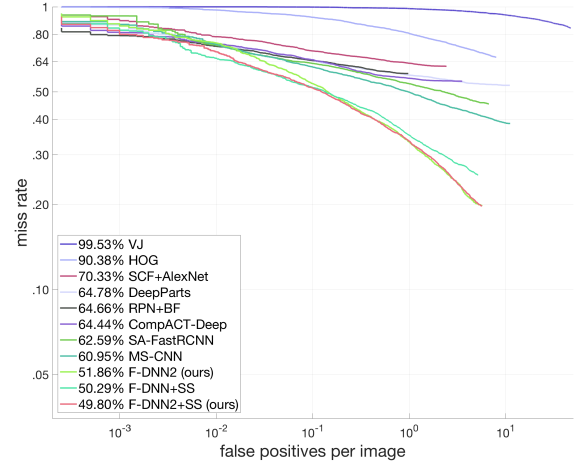


Fig. 6: L-AMR vs. FPPI plot under the 'All' evaluation setting Caltech Pedestrian dataset.

We refer to the new models of this paper as F-DNN2, which is the proposed pedestrian detection system with a fusion network, and F-DNN2+SS, which is F-DNN2 system fused with the semantic segmentation system, and to the models of our previous work [9] as FDNN and FDNN+SS, for fusion of the CG with the classification system only or with both the classification system and the SS network, respectively, as described in Subsections 2.2.2, 2.2.3, and 3.1. Descriptions of each dataset and its evaluation results are given below.

Evaluation on the Caltech Pedestrian data: The Caltech Pedestrian dataset contains 11 sets (S0-S10), where each set consists of 6 to 13 one-minute long videos collected from a vehicle driving through an urban environment. There are about 250,000 frames with about 350,000 annotated BBs and 2,300 unique pedestrians. Each bounding box is assigned with one of the three labels: 'Person', 'People' (large group of individuals), and 'Person?' (unclear identifications). The detailed breakdown performances of our two models (detection system only and detection system plus semantic segmentation system) on this dataset is shown in Table 2. We compare with all the state-of-the-art methods reported on Caltech Pedestrian website. We can see that both of our models significantly outperform others on almost all evaluation settings. On the 'Reasonable' setting, our 'FDNN+SS' best model achieves 8.18% L-AMR, which has a 14.6% relative improvement from the previous best result 9.58% by RPN+BF. Even more accuracy can be obtained by our proposed 'FDNN2' and 'FDNN2+SS' models, which achieve 8.12% and 7.67% L-AMR, respectively. On the 'All' evaluation setting, we achieve 50.29% with 'FDNN+SS', a relative improvement of 17.5% from 60.95% by MS-CNN [19]. The 'FDNN2+SS' has an even better accuracy with an L-AMR of 49.8% under the 'ALL' evaluation setting. The L-AMR vs. FPPI plots for the 'Reasonable' and 'All' evaluation settings are shown in Figure 5 and Figure 6. Except the VJ [22] and the HOG [23] methods, which are plotted as the baselines, all the other results are CNN-based methods.

Evaluation on INRIA: We evaluate the proposed method using the converted INRIA pedestrian dataset provided by Caltech Pedestrian group. There are 614 full positive training images and 288 full positive testing images in the INRIA dataset. At least one

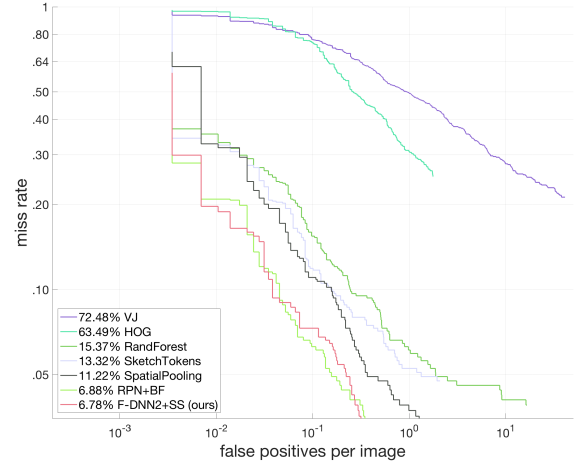


Fig. 7: L-AMR vs. FPPI plot on 'Reasonable' evaluation setting on INRIA dataset.

pedestrian is annotated in each image. To test the generalization capability of our model, we directly test our Caltech-pretrained model on the INRIA test set without any fine-tuning on the INRIA training set. On the 'Reasonable' setting, our 'FDNN2+SS' method achieves 6.78% L-AMR, outperforming the previous best result 6.9% by RPN+BF. Table 3 shows the best results reported on INRIA dataset and Figure 7 shows the L-AMR vs. FPPI plot. Results from VJ and HOG are plotted as the baselines.

Evaluation on ETH: There are 1,804 images from three video sequences in the ETH pedestrian dataset. As we used the ETH images to train our model, in order to test on the ETH dataset, we removed all the training images from the ETH dataset in our training set and retrained our model. On the 'Reasonable' setting, our method achieves 30.02% L-AMR, outperforming the previous best result 30.23% by RPN+BF. Table 4 shows the best results reported on the ETH dataset and Figure 8 shows the L-AMR vs. FPPI plot. Results from VJ and HOG are plotted as the baselines.

Evaluation on KITTI: We further generalize our method to multi-class detection problem and test on KITTI object detection

Method	Reasonable	All	Far	Medium	Near	Occ. none	Occ. partial	Occ. heavy
SCF+AlexNet [3]	23.32%	70.33%	100%	62.34%	10.16%	19.99%	48.47%	74.65%
SAF R-CNN [1]	9.68%	62.6%	100%	51.8%	0%	7.7%	24.8%	64.3%
MS-CNN [19]	9.95%	60.95%	97.23%	49.13%	2.60%	8.15%	19.24%	59.94%
DeepParts [20]	11.89%	64.78%	100%	56.42%	4.78%	10.64%	19.93%	60.42%
CompACT-Deep [21]	11.75%	64.44%	100%	53.23%	3.99%	9.63%	25.14%	65.78%
RPN+BF [6]	9.58%	64.66%	100%	53.93%	2.26%	7.68%	24.23%	69.91%
F-DNN+SS [9]	8.18%	50.29%	77.47%	33.15%	2.82%	6.74%	15.11%	53.76%
F-DNN2 (ours)	8.12%	51.86%	77.99%	36.72%	1.68%	6.75%	17.51%	40.84%
F-DNN2+SS (ours)	7.67%	49.80%	75.83%	35.09%	1.51%	6.35%	16.17%	39.84%

TABLE 2: Detailed breakdown performance comparisons of our models and other state-of-the-art models on the 8 evaluation settings. All numbers are reported in L-AMR.

Method	RPN+BF	SketchTokens	SpatialPooling	RandForest	VJ	HOG	F-DNN2+SS (ours)
L-AMR	6.88%	13.32%	11.22%	15.37%	72.48%	63.49%	6.78%

TABLE 3: Performance comparisons of our models and other state-of-the-art models on the INRIA dataset.

Method	RPN+BF	TA-CNN	SpatialPooling	RandForest	VJ	HOG	F-DNN2+SS (ours)
L-AMR	30.32%	34.98%	37.37%	45.04%	74.69%	89.89%	30.02%

TABLE 4: Performance comparisons of our models and other state-of-the-art models on the ETH dataset.

Setting	Description
Easy	Min. BB height: 40 Px, Max. occlusion level: Fully visible, Max. truncation: 15%
Moderate	Min. BB height: 25 Px, Max. occlusion level: Partly occluded, Max. truncation: 30%
Hard	Min. BB height: 25 Px, Max. occlusion level: Difficult to see, Max. truncation: 50%

TABLE 5: Evaluation settings for KITTI object detection dataset.

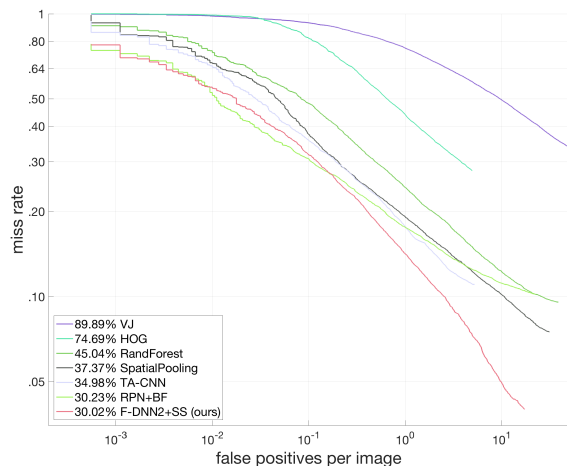


Fig. 8: L-AMR vs. FPPI plot on 'Reasonable' evaluation setting on ETH dataset.

dataset. KITTI object detection dataset contains 7,481 training images and 7,518 test images. All the annotations are split into 7 classes such as cars, vans, trucks, pedestrians, cyclists, trams, and 'Don't care'. Only cars, pedestrians, and cyclists are evaluated. There are three evaluation settings as shown in Table 5. Following [19], we split the training data into a training set and a validation set. We fine-tune three models using all training annotations for the three evaluating classes respectively. For the three models, we set the main aspect ratio to the mean aspect ratio of each class, which is 0.4 for pedestrians, 0.7 for cyclists, and 1.6 for cars. Table 6 shows the results on KITTI object detection dataset.

We achieved comparable results on all classes. Since the Caltech Pedestrian dataset doesn't distinguish between pedestrians and cyclists, while the KITTI object detection dataset does, it degrades our performance on the pedestrians class and the cyclists class on KITTI.

Benchmark	Easy	Moderate	Hard
Car	89.68 %	85.11 %	70.35 %
Pedestrian	74.05 %	61.17 %	57.15 %
Cyclist	67.06 %	51.85 %	46.67 %

TABLE 6: Evaluation results on KITTI object detection dataset.

Method	Reasonable
CG	13.06%
CG+GglNet	8.64%
CG+ResNet	8.38%
CG+GglNet+ResNet+Fusion net (F-DNN2)	8.12%
CG+GglNet+ResNet+Fusion net+SS (F-DNN2+SS)	7.67%

TABLE 7: Ablation study of our system.

Method	Hard-label	Soft-label
CG + ResNet	8.97%	8.38%
CG + GglNet	9.41%	8.64%
CG + GglNet + ResNet	8.65%	8.12%

TABLE 8: Effectiveness of the soft-label method compared to the conventional hard-label method on Caltech Pedestrian dataset using the reasonable setting.

Method	Speed on TITAN X (seconds per image)
CompACT-Deep	0.5
SAF R-CNN	0.59
F-DNN2	0.3
F-DNN2 (Reasonable)	0.16
CG+GglNet (Reasonable)	0.11
CG+SqueezeNet (Reasonable)	0.09
F-DNN2+SS	2.48

TABLE 9: A comparison of speed among the state-of-the-art models.

4.3 Result analysis

4.3.1 An ablation study: effectiveness of the network fusion technique

In this subsection, we analysis the performance increases step by step from the candidate generator (CG) to the final system. The L-AMR is 13.06% by using the candidate generator alone, due to the large number of false positives. By fusing the candidate generator with GoogLeNet, we can improve the performance to 8.64%. By fusing the candidate generator with ResNet-50, we can improve the L-AMR to 8.38%. Furthermore, by fusing the candidate generator with GoogLeNet and ResNet-50 using our proposed fusion net, we can achieve the lowest L-AMR so far at 8.12%. Finally, by fusing the semantic segmentation network into our system, we can achieve the best performance at 7.67%. The analysis shows the capability of the network fusion framework and the advantages of using the idea of ensemble learning. The results of ablation study are given in Table 7.

4.3.2 GoogLeNet VS. ResNet-50

We explore how each part of the classification system contributes to our final results. The breakdown performance comparisons of all evaluation settings on the Caltech Pedestrian dataset between fusing with GoogLeNet alone and fusing with ResNet-50 alone are given in Table 10. From the results we can see that GoogLeNet works better in partial/heavy occluded pedestrians while ResNet-50 works better in non-occluded pedestrians. By analyzing the weights learnt in Equation (10), we see that the weight for GoogLeNet is 1.11 and the weight for ResNet-50 is 2.22, which means that our fusion network values the ResNet-50 more than the GoogLeNet. This is reasonable since there are more non-occluded pedestrians in the training data.

4.3.3 Soft-label method versus hard-label method

To test how effectively the soft-label method improves the performance, we compare with the conventional hard-label method on the Caltech Pedestrian dataset. Since we use the overlap ratio between the candidate bounding box and the ground-truth annotation to assign labels, the soft-label method gives us not only the information of the existence of a pedestrian in each candidate’s bounding box, but also how much of the bounding box belongs to the pedestrian. This feature benefits even more in cases where the overlap ratio is around 0.5: e.g. it is too risky to directly assign a hard-label 1 or 0 to a bounding box with overlap ratio 0.49 or 0.51. We give the performance comparisons between the hard-label method and the soft-label method in Table 8.

4.3.4 Results visualization on challenging scenarios and failure cases

We visualize the detection results generated by our system compared with previous state-of-the-art systems on several challenging

scenarios: far pedestrians, crowded scene, occluded pedestrians, pedestrians overlap with each other. Figure 9 visualizes detection results on five challenging scenarios. In Figure 9, each row represents one challenging scenario and the four columns show the bounding boxes from the ground truth, RPN+BF, F-DNN, and F-DNN2, respectively. From the visualizations we can see that our model is more robust and accurate on various challenging cases.

4.3.5 Speed analysis

We use one NVIDIA TITAN X GPU to analysis the processing speed of each component and the overall architecture of our network. There are four main components: the candidate generator, GoogLeNet, ResNet-50, and the semantic segmentation network. Since the candidate generator has a fully convolutional framework which removes the fully connected layers in the original VGG net, it takes 0.06s to process one image. To test the processing time of the classification system, we have two settings: the first test runs on all pedestrian candidates; The second test runs only on candidates above 40-pixel in height, which is designed for the ‘reasonable’ evaluation setting. Since the number of pedestrian candidates varies from image to image, the test reports the mean processing time of all images. Running the candidate generator followed by GoogLeNet and ResNet-50 in parallel on one GPU, the speed for the classification system is 0.3s and 0.16s per image for the two tests. To achieve real time pedestrian detection, we fuse the candidate generator with SqueezeNet [24]. The processing time per image is 0.09s, while being 10.8% in L-AMR. By processing the semantic segmentation network in parallel with the pedestrian detection system and fusing them together, the processing time per image of our final model is 2.48s. The speed comparisons of our models and other methods are given in Table 9.

5 CONCLUSION AND DISCUSSION

We present an effective solution to the pedestrian detection problem in this paper. The proposed system consists of two parallel subsystems: The main pedestrian detection system to generate all detections and the semantic segmentation system to help refine the results. The pedestrian detection system further consists of a pedestrian candidate generator and a classification system. To give more bounding box information to the classification network, we propose a new soft-label method which assigns floating point label to all classes. We implement the idea of ensemble learning to design the classification system. We proposed a soft-rejection network fusion methodology to combine the opinions of all classification networks and the semantic segmentation network with that of the candidate generator network.

The performance of our system is evaluated on four popular pedestrian detection datasets: Caltech Pedestrian dataset, INRIA dataset, ETH dataset, and KITTI dataset. We achieve the state-of-the-art performance on the first three datasets and comparable results on the KITTI dataset. Our system also works faster in processing speed than other state-of-the-arts when testing using the same experiment settings. The results and analysis show that our system works accurately, efficiently, and robustly to detect pedestrians and other object classes under various challenging scenarios.

REFERENCES

- [1] J. Li, X. Liang, S. Shen, T. Xu, and S. Yan, “Scale-aware fast R-CNN for pedestrian detection,” *CoRR*, vol. abs/1510.08160, 2015. [Online]. Available: <http://arxiv.org/abs/1510.08160>

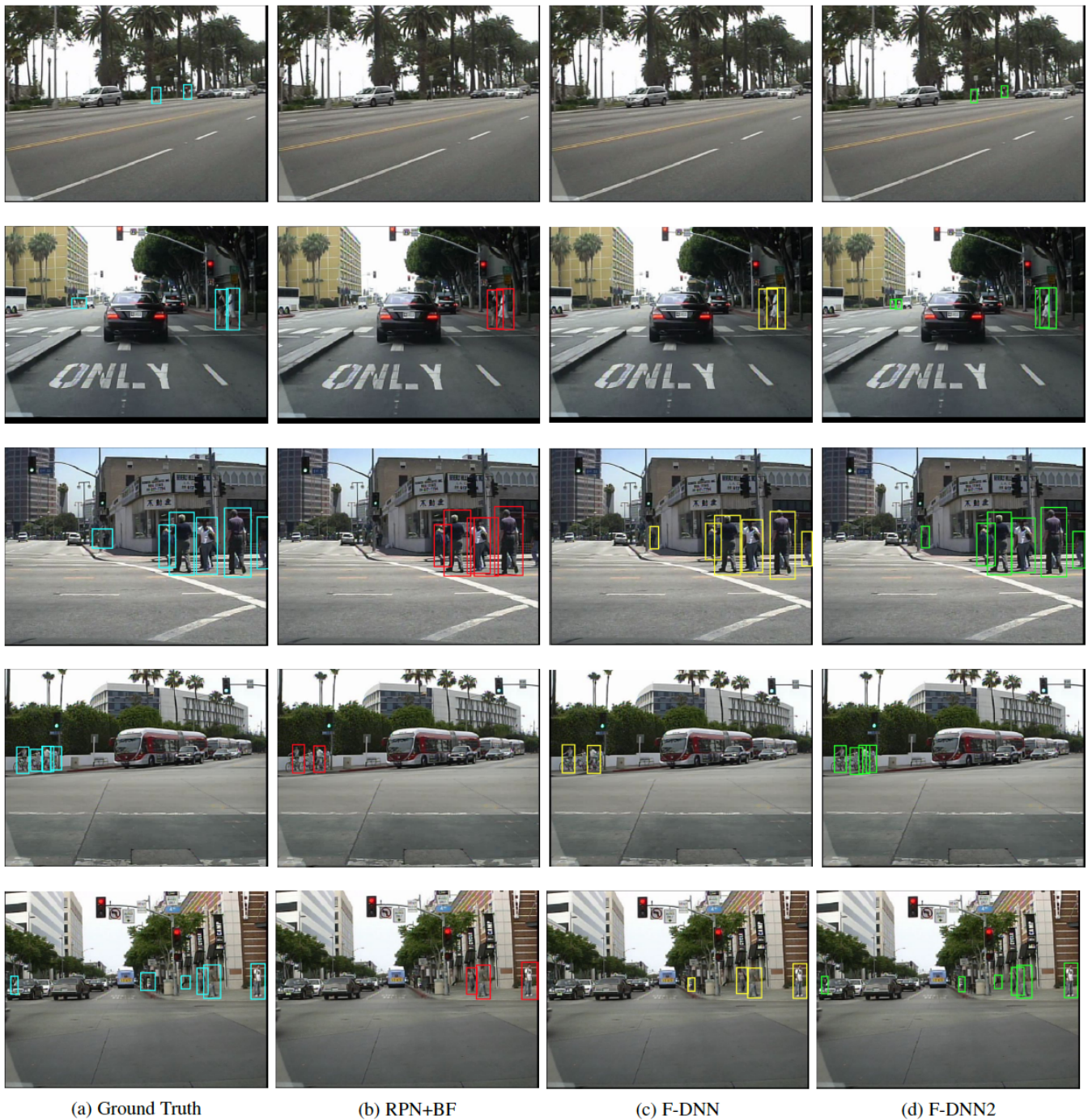


Fig. 9: Detection comparisons on five challenging pedestrian detection scenarios. Each row represents one challenging scenario and the four columns show the bounding boxes from the ground truth, RPN+BF, F-DNN, and F-DNN2, respectively.

Method	Reasonable	All	Far	Medium	Near	Occ. none	Occ. partial	Occ. heavy
CG+GglNet	8.64%	51.59%	76.87%	37.75%	1.72%	7.18%	18.05%	41.19%
CG+ResNet	8.38%	49.58%	74.60%	34.88%	1.70%	6.94%	19.26%	42.71%

TABLE 10: Breakdown comparisons between CG+GglNet and CG+ResNet on Caltech Pedestrian dataset.

- [2] R. Girshick, "Fast R-CNN," in *International Conference on Computer Vision (ICCV)*, 2015.
- [3] J. H. Hosang, M. Omran, R. Benenson, and B. Schiele, "Taking a deeper look at pedestrians," *CoRR*, vol. abs/1501.05790, 2015. [Online]. Available: <http://arxiv.org/abs/1501.05790>
- [4] R. Benenson, M. Omran, J. H. Hosang, and B. Schiele, "Ten years of pedestrian detection, what have we learned?" *CoRR*, vol. abs/1411.4304, 2014. [Online]. Available: <http://arxiv.org/abs/1411.4304>
- [5] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds. Curran Associates, Inc., 2012, pp. 1097–1105. [Online]. Available: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>
- [6] L. Zhang, L. Lin, X. Liang, and K. He, "Is faster R-CNN doing well for pedestrian detection?" *To appear in ECCV 2016*, 2016.
- [7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems (NIPS)*, 2015.
- [8] P. Dollár, "Quickly boosting decision trees - pruning underachieving features early," in *ICML. International Conference on Machine Learning*, June 2013. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/quickly-boosting-decision-trees-pruning-underachieving-features-early/>
- [9] X. Du, M. El-Khamy, J. Lee, and L. S. Davis, "Fused DNN: A deep neural network fusion approach to fast and robust pedestrian detection," *CoRR*, vol. abs/1610.03466, 2016. [Online]. Available: <http://arxiv.org/abs/1610.03466>
- [10] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single shot multibox detector," *arXiv:1512.02325*, 2015.
- [11] W. Liu, A. Rabinovich, and A. C. Berg, "ParseNet: Looking wider to see better," *CoRR*, vol. abs/1506.04579, 2015.
- [12] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *To appear in ICLR 2016*, 2016.
- [13] 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," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [14] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: common objects in context," *CoRR*, vol. abs/1405.0312, 2014. [Online]. Available: <http://arxiv.org/abs/1405.0312>
- [15] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.
- [16] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. E. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," *CoRR*, vol. abs/1409.4842, 2014. [Online]. Available: <http://arxiv.org/abs/1409.4842>
- [17] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," *arXiv preprint arXiv:1408.5093*, 2014.
- [18] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: An evaluation of the state of the art," *PAMI*, vol. 34, 2012.
- [19] Z. Cai, Q. Fan, R. Feris, and N. Vasconcelos, "A unified multi-scale deep convolutional neural network for fast object detection," in *ECCV*, 2016.
- [20] Y. Tian, P. Luo, X. Wang, and X. Tang, "Deep learning strong parts for pedestrian detection," in *ICCV*, 2015.
- [21] Z. Cai, M. Saberian, and N. Vasconcelos, "Learning complexity-aware cascades for deep pedestrian detection," in *ICCV*, 2015.
- [22] P. Viola and M. J. Jones, "Robust real-time face detection," *Int. J. Comput. Vision*, vol. 57, no. 2, pp. 137–154, May 2004. [Online]. Available: <http://dx.doi.org/10.1023/B:VISI.0000013087.49260.fb>
- [23] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *In CVPR*, 2005, pp. 886–893.
- [24] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer, "Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size," *arXiv:1602.07360*, 2016.