# NMS-Loss: Learning with Non-Maximum Suppression for Crowded Pedestrian Detection

Zekun Luo
Youtu Lab, Tencent
Shanghai, China

Zheng Fang
Beihang University
Beijing, China

Sixiao Zheng
Fudan University
Shanghai, China

Yabiao Wang
Youtu Lab, Tencent
Shanghai, China

Yanwei Fu*
Fudan University
Shanghai, China

## ABSTRACT

Non-Maximum Suppression (NMS) is essential for object detection and affects the evaluation results by incorporating False Positives (FP) and False Negatives (FN), especially in crowd occlusion scenes. In this paper, we raise the problem of weak connection between the training targets and the evaluation metrics caused by NMS and propose a novel NMS-Loss making the NMS procedure can be trained end-to-end without any additional network parameters. Our NMS-Loss punishes two cases when FP is not suppressed and FN is wrongly eliminated by NMS. Specifically, we propose a pull loss to pull predictions with the same target close to each other, and a push loss to push predictions with different targets away from each other. Experimental results show that with the help of NMS-Loss, our detector, namely NMS-Ped, achieves impressive results with Miss Rate of 5.92% on Caltech dataset and 10.08% on CityPersons dataset, which are both better than state-of-the-art competitors.

## CCS CONCEPTS

• **Computing methodologies** → Object detection.

## KEYWORDS

pedestrian detection, loss function, Non-Maximum suppression

*Yanwei Fu is with the School of Data Science and MOE Frontiers Center for Brain Science, Fudan University, Shanghai 200433, China (e-mail: yanweifu@fudan.edu.cn).

## 1 INTRODUCTION

Pedestrian detection [12] is an essential computer vision task that has numerous applications such as automatic driving, video surveillance and person re-identification. With the help of deep convolution neural networks (CNNs) [17, 29, 37], the performance of pedestrian detection has been significantly improved. However, the False Negatives (FN) in crowd occlusion scenes and False Positives (FP) generated for the same person are still the fundamental challenges.

Existing methods for pedestrian detection can mainly be divided into two categories: hand-crafted feature based [10, 11, 15, 16, 24, 33, 34, 37] and deep learning based [3, 4, 14, 20, 23, 27, 31, 35]. The first one applies the sliding-window way to get different scales of patches, then uses human-designed feature extractor such as Haar [30] and HoG [9] to obtain feature representation, last utilizes SVM [8] classifier to filter background. These hand-crafted feature representations could not handle complex scenes. The second one uses deep convolutional neural networks (CNNs) to obtain high-level semantic feature representation, which has a discriminative ability to deal with complex scenes for pedestrian detection. To alleviate FN issue in high occlusion scenes, different variants of Non-Maximum Suppression (NMS) [1, 18, 19] are proposed to change NMS threshold during inference adaptively. To reduce FP, many works [5, 6] jointly predict pedestrian boxes and parts information such as head due to that it is less occluded. However, the objective between training and inference is inconsistent, which may result in sub-optimal performance for pedestrian detection.

NMS is an essential procedure for object detection tasks. Modern pedestrian detectors rely on NMS to remove duplicate detections for both one-stage and two-stage approaches. The nearby detections around one object will be removed once its interaction over union (IoU) with the object is larger than the pre-defined threshold. During the training process, there is no such process, thus resulting in inconsistency between optimized detection training results and final inference results. To handle the inconsistency problem, NMS process should be incorporated into the training process. To this end, we propose a novel NMS-Loss. There are two components, pull and push losses, in our NMS-Loss. Pull loss aims to raise the precision by pulling FP close to the max score prediction, and push loss focuses on improving recall by pushing predictions away from each other. With the help of NMS-Loss, false predictions on the evaluation metric can be directly reflected on loss functions, and thus be directly optimized.

The main contribution of this work lies in the following aspects.

- We firstly raise the problem of weak connection between training targets and evaluation metrics in pedestrian detection and propose a novel NMS-Loss making the NMS procedure can be trained end-to-end, which does not introduce any parameters nor runtime cost.
- We propose finely designed pull and push losses helping the network to boost performance on precision and recall, respectively, which considering both prediction coordinates and confidence.
- With the help of NMS-Loss, in pedestrian detection, our proposed NMS-Ped outperforms SOTA methods on the widely used Caltech and CityPersons datasets.

## 2 NMS-LOSS

### 2.1 Overview of NMS-Loss

The traditional NMS procedure is shown in Alg. 1 without considering the red texts. Starting with a set of detection boxes $\mathcal{B}$ with corresponding scores $\mathcal{S}$, NMS firstly moves the proposal $b_m$ with the maximum score from the set $\mathcal{B}$ to the set of final kept detections $\mathcal{K}$. It then removes any box in $\mathcal{B}$ and its score in $\mathcal{S}$ that has an overlap with the $b_m$ higher than a manually set threshold $N_t$. This process is repeated for the remaining $\mathcal{B}$ set.

However, no existing approaches take the NMS into the training process to adjust the detection boxes, making the learning targets inconsistent with the evaluation metric, which means FP not suppressed by NMS and FN eliminated by NMS can harm the precision and recall, respectively. To avoid inconsistency, we propose the NMS-Loss taking the NMS procedure into the training process, which adaptively selects the false predictions caused by NMS and uses two well-designed *pull* and *push losses* to minimize the FP and FN, respectively. Specifically, our NMS-Loss is defined as:

$$L_{nms} = \lambda_{pull}L_{pull} + \lambda_{push}L_{push}, \quad (1)$$

where $L_{pull}$ is the pull loss to punish the FP not suppressed by NMS and $L_{push}$ is the push loss to punish the FN wrongly eliminated by NMS. Coefficients $\lambda_{pull}$ and $\lambda_{push}$ are the weights for balancing losses. Details of our NMS-Loss are present in Algorithm 1 emphasized with red color. Different from the traditional NMS, we use a set $\mathcal{G}$ containing corresponding ground truth indexes of detection boxes, which is used to identify FP and FN. In the NMS-Loss calculating procedure, $\mathcal{M}$ is an auxiliary dictionary with the ground truth index as key and corresponding max score detection as value, which is used to record the max score prediction of each ground truth. Our NMS-Loss is naturally merged into the NMS procedure without incorporating any additional training parameters. The runtime cost of NMS-Loss is zero for testing.

### 2.2 Pull Loss Definition

With the objective to reduce FP, we need to find out wrongly kept predictions. To this end, in every iteration, we check whether the current max score prediction $b_m$ is the max score prediction for its corresponding $g_m$ ground truth. If not, it means $b_m$ is an FP not suppressed by NMS, pull loss should be performed between $b_m$ and the max score prediction $b_{max}$ of the $g_m$ ground truth (see Fig. 1). Formally, our pull loss is calculated as:

$$L_{pull} = -ln(1 - N_t + IoU(b_{max}, b_m))s_m, \quad (2)$$

**Algorithm 1: NMS-Loss Calculating Procedure**

**Input:**
$\mathcal{B} = [b_1, \ldots, b_N], \mathcal{S} = [s_1, \ldots, s_N], N_t, \mathcal{G} = [g_1, \ldots, g_N]$
$\mathcal{B}$ is the list of initial detection boxes
$\mathcal{S}$ contains corresponding detection scores
$N_t$ is the NMS threshold
$\mathcal{G}$ contains corresponding ground truth indexes

**Auxiliary Variable:**
$\mathcal{K} \leftarrow [\ ], \mathcal{M} \leftarrow dictionary(),$
$\mathcal{K}$ is the list to keep final detections after NMS
$\mathcal{M}$ is a dictionary using the ground truth index as key and corresponding max score detection as value

**begin**
  **while** $\mathcal{B} \neq empty$ **do**
    $m \leftarrow argmax\ \mathcal{S}$ ;
    **if** $g_m$ *not in* $\mathcal{M}.keys()$ **then**
      $\mathcal{M}[g_m] \leftarrow b_m$;
    **else**
      $b_{max} \leftarrow \mathcal{M}[g_m]$;
      pull_loss($b_{max}, b_m$);   Eq. (2)
    **end**
    $\mathcal{K} \leftarrow \mathcal{K} \cup b_m$ ; $\mathcal{B} \leftarrow \mathcal{B} - b_m$;
    $\mathcal{S} \leftarrow \mathcal{S} - s_m$ ; $\mathcal{G} \leftarrow \mathcal{G} - g_m$;
    **for** $b_i$ *in* $\mathcal{B}$ **do**
      **if** $IoU(b_m, b_i) \geq N_t$ **then**
        **if** $g_m \neq g_i$ **then**
          push_loss($b_m, b_i$);   Eq. (3)
        **end**
        $\mathcal{B} \leftarrow \mathcal{B} - b_i$; $\mathcal{S} \leftarrow \mathcal{S} - s_i$; $\mathcal{G} \leftarrow \mathcal{G} - g_i$;
      **end**
    **end**
  **end**
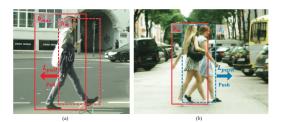  **return** $\mathcal{K}$
**end**



**Figure 1: Illustration of our NMS-Loss. All boxes $b_{max}$, $b_m$ and $b_i$ are predictions as described in Alg. 1, where boxes with the same color have the same target and boxes with the solid line get a higher score than boxes with the dotted line. In (a), $b_m$ is a FP not suppressed by $b_{max}$. Our $L_{pull}$ pulls $b_m$ towards $b_{max}$. In (b), $b_i$ is a FN wrongly eliminated by $b_m$. Our $L_{push}$ pushes $b_i$ away from $b_m$.**

where $N_t$ is the predefined NMS threshold and $s_m$ is the prediction score corresponding to $b_m$. We note two properties of the pull loss: (1) When the IoU between $b_{max}$ and $b_m$ is small, pull loss tends to increase, forcing the network to learn to pull $b_m$ toward $b_{max}$. The NMS threshold $N_t$ is used to prevent the gradient of outliers influence too much on model learning. Besides, for the NMS procedure, we just need to make the IoU between FP and TP higher than $N_t$. Using $N_t$ in pull loss to reduce the gradient of

outliers can make the network easy to learn. (2) The prediction score of FP can also have a strong effect on pull loss. FP with a higher score has a greater impact on evaluation results and intuitively needs to be paid more attention. Besides, it makes the network learn to fix FP not only just conditioning the box coordinates but also considering lower the prediction scores.

## 2.3 Push Loss Definition

In NMS, the current max score prediction $b_m$ eliminates boxes which get an IoU higher than $N_t$ with $b_m$. If the eliminated box $b_i$ corresponds to different ground truth index with $b_m$, $b_i$ will be a FN and reduce recall (see Fig. 1). To avoid $b_i$ from being wrongly eliminated, we propose a push loss to penalize FN:

$$L_{push} = -ln(1 - IoU(b_i, b_m))s_i, \qquad (3)$$

where $s_i$ is the prediction score corresponding to $b_i$. Different from pull loss, as $IoU(b_i, b_m) \rightarrow 1$, the push loss goes higher and the model learns to push $b_i$ away from $b_m$. To avoid the model tending to reduce the push loss by lowering the score of FN, we use the $s_i$ only for reweighting losses without back propagating gradient.

For crowded scenes, especially in the CityPersons dataset, the ground truths of bounding boxes are overlapped with each other. It is unreasonable to push their predictions away from each other with an IoU equals to zero. To handle this problem, we only calculate $L_{push}$ on prediction whose IoU is higher than the IoU of its corresponding ground truth boxes.

Our pull and push loss are performed on predictions. When the pull/push loss is activated, the network tries to pull/push both predictions close to/away from each other, respectively. Since high score predictions generally get a more accurate location, it is unreasonable to move an accurate prediction based on an inaccurate one. To handle this, we stop the gradient backward propagation of high score predictions, leading the network to focus on false predictions.

## 3 EXPERIMENTS

### 3.1 Experimental Setup

**Datasets and Evaluation metrics.** We evaluate our method on two challenging pedestrian datasets: Caltech [12, 13] and CityPersons [36]. We report performance using standard average-log MR between $[10^{-2}, 10^0]$ of False Positive per Image (FPPI). A minimum IoU threshold of 0.5 is required for detected box to match with a ground truth box. By default, we report the results on Reasonable subsets is a widely used setup where the pedestrian is at least $65\%$ visible and 50 pixels tall.

**Experimental Settings.** As shown in RPN+BF [35], small instances are hard to be detected in the low-resolution feature maps provided by RoI-Pooling, which is more severe in pedestrian detection. Therefore, we used Faster R-CNN [28] as our baseline, but made two adjustments: (1) Inspired by [35], we use a separate network to construct the RCNN and put the cropped original image to RCNN for further refinement. This improves the ability of the network to detect small instances, but it is not suitable for instances with large scale changes. (2) There is an additional weak semantic segmentation loss [3] to boost performance. Note that the baseline has the same settings as our NMS-Ped except that there is no NMS-Loss in baseline.
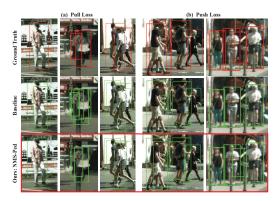


**Figure 2: Comparison between the cases with/without using pull/push loss. Green bounding boxes are predicted pedestrians whose score is greater than** 0.8 **and red bounding boxes are ground truth. Our pull loss effectively suppresses FP in both sparse scenes and crowded scenes (left three columns), yielding higher precision. Our push loss robustly handles occlusions (right two columns), yielding higher recall.**

**Table 1: Comparison of our NMS-Ped with the baseline on CityPersons.**

| Method | MR |
|---|---|
| baseline | 11.20% |
| baseline + pull loss | 10.58% |
| baseline + push loss | 10.61% |
| **NMS-Ped** | **10.08%** |

**Table 2: Comparison on different thresholds $N_t$ of NMS-Loss on CityPersons.**

| $N_t$ | 0.4 | 0.45 | **0.5** | 0.55 |
|---|---|---|---|---|
| MR | 10.76% | 10.66% | **10.08%** | 10.67% |

PyTorch [26] is used to train the NMS-ped for both datasets. We use 8 NVIDIA GPUs with a mini-batch comprises 1 image per GPU. SGD with momentum of 0.9 and weight decay of $1 \times 10^{-4}$ is adopted for training. Both datasets are trained only using the images with foreground. Random cropping and flipping are used for data augmentation. Detailed settings on Caltech and CityPersons are described as follows:

**Caltech:** The learning rate for Caltech is $5 \times 10^{-3}$ and is dropped by a factor of 10 after $9,600$ iterations and $13,200$ iterations. The images are resized to $1280 \times 960$ in our experiments. The weights for pull loss and push loss are both 0.1 getting from experiments.

**CityPersons:** The learning rate for CityPersons is $1 \times 10^{-2}$ and dropped by a factor of 10 after $24,000$ iterations and $33,000$ iterations. We use the original image resolution of $2048 \times 1024$ in our experiments. The weights for pull and push loss are 0.1 and 0.001 respectively for the reason that CityPersons contains much more crowded scenes than Caltech and lots of instances are heavily overlapped with others. Giving a relatively lower weight for push loss will reduce the gradient of pushing and make multi-tasks work well.

**Table 3: Comparison on CityPersons dataset.**

| Method | Backbone | MR |
|---|---|---|
| RepLoss [31] | ResNet-50 | 13.20% |
| OR-CNN [38] | ResNet-50 | 12.80% |
| Adaptive-NMS [19] | VGG-16 | 11.90% |
| CSP [21] | ResNet-50 | 11.00% |
| MGAN [25] | VGG-16 | 11.50% |
| $R^2$NMS [18] | VGG-16 | 11.10% |
| EMD-RCNN [7] | ResNet-50 | 10.70% |
| Our baseline | ResNet-50 | 11.20% |
| **NMS-Ped** | **ResNet-50** | **10.08%** |

**Table 4: Comparisons on Caltech dataset.**

| Method | Backbone | MR |
|---|---|---|
| RPN+BF [35] | VGG-16 | 9.58% |
| F-DNN [14] | ResNet-50 | 8.65% |
| SDS-RCNN [3] | VGG-16 | 7.36% |
| MGAN [25] | VGG-16 | 6.83% |
| AR-Ped [2] | VGG-16 | 6.45% |
| SSA-CNN [39] | VGG-16 | 6.27% |
| TFAN+TDEM+PRM [32] | ResNet-101 | 6.50% |
| $W^2$Net [22] | ResNet-50 | 6.37% |
| Our baseline | ResNet-50 | 6.61% |
| **NMS-Ped** | **ResNet-50** | **5.92%** |

## 3.2 Ablation Studies

We conduct experiments on CityPersons to evaluate our NMS-Loss for the reason that pedestrian in CityPersons is more crowded and challenging. There are enough complicated scenes to review effectiveness of our approach.

**Baseline comparison.** Tab. 1 shows the performance of our baseline with separate components. When only the pull loss is used, MR can be reduced from 11.20% to 10.58%. Fig. 2 shows some results corrected for using pull loss. In both sparse scenes (first column) and crowded scenes (second and third columns), our pull loss will effectively pull predictions targeting on the same ground truth close to each other. The same experiments are conducted on push loss. With the help of push loss, the MR can be reduced from 11.20% to 10.61%. Some visible results are present in Fig. 2 showing the corrected predictions for using push loss. In the occlusion scenes (right two columns), push loss trained model performs more robust, even detected the unlabeled instance (fourth column). When we use the complete NMS-Loss, our NMS-Ped can be boosted from both pull loss and push loss, getting an amazing 10.08% MR.

**Experiments on hyperparameters.** Tab. 2 shows our results with different thresholds $N_t$ on NMS-Loss. When $N_t$ is lower than evaluation metric threshold 0.5, push loss will be activated more frequently and pull loss will not be activated making the network produce more FPs that harms precision. In contrast, when $N_t$ is higher than 0.5, more FNs will be produced and lower recall. Our NMS-Loss performs robust with various NMS thresholds, gaining stable improvement. When we use $N_t$ equivalent to the threshold 0.5, our NMS-Loss yields the best performance.

**Table 5: Comparison between RepLoss and NMS-Loss on the CityPersons. We use $MR_b$, $MR$, $MR_i$, $MR_r$ to represent the $MR$ of baseline model, $MR$ of complete model, $MR$ of the improvement and relative improvement based on the baseline, respectively.**

| Method | Backbone | $MR_b \downarrow$ | $MR \downarrow$ | $MR_i \uparrow$ | $MR_r \uparrow$ |
|---|---|---|---|---|---|
| RepLoss | ResNet-50 | 14.6% | 13.2% | 1.4% | 9.59% |
| NMS-Ped | ResNet-50 | 11.2% | 10.08% | 1.12% | **10.00%** |

## 3.3 Comparisons with SOTA methods

To demonstrate the effectiveness of our NMS-Loss, we compare NMS-Ped with the SOTA methods on CityPersons and Caltech. Tab. 3 presents the performance of NMS-Ped and SOTA methods on the CityPersons dataset. With the help of NMS-Loss, our method improve the MR of baseline from 11.20% to 10.08%, better than the SOTA method EMD-RCNN [7] (MR of 10.70%). Tab. 4 presents the performance on Caltech, the MR of NMS-Ped is 5.92%, better than SOTA method $W^2$Net [22] (MR of 6.37%). With the help of NMS-Loss, we can obtain more than 10% improvement in NMS-Ped compared with baseline. This demonstrates the effectiveness of our NMS-Loss.

## 3.4 Difference to RepLoss

We make a detailed comparison between our NMS-Loss and the RepLoss [31] for the reason that both methods pull and push predictions based on their targets. There are three main differences: (1) RepLoss is performed on all instances, while NMS-Loss is only performed on instances wrongly processed by NMS, which enables end-to-end training. (2) RepLoss only considers regression, while the score is also used in NMS-Loss to reweight instances. (3) In dense crowd scenarios, RepLoss pushes instances away even if their targets are originally close to each other, making the repulsion loss contradicts with the regression loss. Instead, NMS-Loss pushes instances whose IoU with others is higher than the IoU of its corresponding ground truth boxes, which can eliminates the contradiction of RepLoss. As shown in Tab. 5, our NMS-Loss not only performs better than RepLoss, but also gains higher relative improvement on CityPersons. This demonstrates that our NMS-Loss can achieve stable relative improvement (higher than 10%) on the widely used datasets.

## 4 CONCLUSION AND FUTURE WORK

In this work, we raise the problem of weak connection between training targets and evaluation metrics in the object detection. To address this, we propose the NMS-Loss which contains two components called pull loss and push loss, making the false predictions can be directly reflected on loss functions. With the help of NMS-Loss, the model can be trained with NMS end-to-end and pay more attention to the false predictions caused by NMS. Our NMS-Loss can be easily incorporated into network, which does not introduce any parameters nor runtime cost. NMS-Loss is only suitable for single class object detection, in the future, we will extend our NMS-Loss to other tasks by further considering object classes in generic detections.

# REFERENCES

[1] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. 2017. Soft-NMS–Improving Object Detection With One Line of Code. In *ICCV*. 5561–5569.

[2] Garrick Brazil and Xiaoming Liu. 2019. Pedestrian Detection with Autoregressive Network Phases. In *CVPR*. 7231–7240.

[3] Garrick Brazil, Xi Yin, and Xiaoming Liu. 2017. Illuminating pedestrians via simultaneous detection & segmentation. In *ICCV*. 4950–4959.

[4] Zhaowei Cai, Quanfu Fan, Rogerio S Feris, and Nuno Vasconcelos. 2016. A unified multi-scale deep convolutional neural network for fast object detection. In *ECCV*. Springer, 354–370.

[5] Cheng Chi, Shifeng Zhang, Junliang Xing, Zhen Lei, Stan Z Li, and Xudong Zou. 2020. Relational learning for joint head and human detection. In *AAAI*. 10647–10654.

[6] Cheng Chi, Shifeng Zhang, Junliang Xing, Zhen Lei, Stan Z Li, Xudong Zou, et al. 2020. PedHunter: Occlusion Robust Pedestrian Detector in Crowded Scenes.. In *AAAI*. 10639–10646.

[7] Xuangeng Chu, Anlin Zheng, Xiangyu Zhang, and Jian Sun. 2020. Detection in Crowded Scenes: One Proposal, Multiple Predictions. In *CVPR*. 12214–12223.

[8] Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine learning* 20, 3 (1995), 273–297.

[9] Navneet Dalal and Bill Triggs. 2005. Histograms of oriented gradients for human detection. In *CVPR*. 886–893.

[10] Piotr Dollár, Ron Appel, Serge Belongie, and Pietro Perona. 2014. Fast feature pyramids for object detection. *PAMI* 36, 8 (2014), 1532–1545.

[11] Piotr Dollár, Zhuowen Tu, Pietro Perona, and Serge Belongie. 2009. Integral channel features. In *BMVC*.

[12] P Dollar, C Wojek, B Schiele, and P Perona. 2009. Pedestrian detection: A benchmark. In *CVPR*. 304–311.

[13] Piotr Dollar, Christian Wojek, Bernt Schiele, and Pietro Perona. 2011. Pedestrian detection: An evaluation of the state of the art. *PAMI* 34, 4 (2011), 743–761.

[14] Xianzhi Du, Mostafa El-Khamy, Jungwon Lee, and Larry Davis. 2017. Fused DNN: A deep neural network fusion approach to fast and robust pedestrian detection. In *WACV*. IEEE, 953–961.

[15] Pedro F Felzenszwalb, Ross B Girshick, and David McAllester. 2010. Cascade object detection with deformable part models. In *CVPR*. IEEE, 2241–2248.

[16] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. 2009. Object detection with discriminatively trained part-based models. *PAMI* 32, 9 (2009), 1627–1645.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*. 770–778.

[18] Xin Huang, Zheng Ge, Zequn Jie, and Osamu Yoshie. 2020. NMS by Representative Region: Towards Crowded Pedestrian Detection by Proposal Pairing. In *CVPR*. 10750–10759.

[19] Songtao Liu, Di Huang, and Yunhong Wang. 2019. Adaptive NMS: Refining Pedestrian Detection in a Crowd. In *CVPR*. 6459–6468.

[20] Wei Liu, Shengcai Liao, Weidong Hu, Xuezhi Liang, and Xiao Chen. 2018. Learning efficient single-stage pedestrian detectors by asymptotic localization fitting. In *ECCV*. 618–634.

[21] Wei Liu, Shengcai Liao, Weiqiang Ren, Weidong Hu, and Yinan Yu. 2019. High-level Semantic Feature Detection: A New Perspective for Pedestrian Detection. In *CVPR*. 5187–5196.

[22] Yan Luo, Chongyang Zhang, Muming Zhao, Hao Zhou, and Jun Sun. 2020. Where, What, Whether: Multi-Modal Learning Meets Pedestrian Detection. In *CVPR*. 14065–14073.

[23] Jiayuan Mao, Tete Xiao, Yuning Jiang, and Zhimin Cao. 2017. What can help pedestrian detection?. In *CVPR*. 3127–3136.

[24] Woonhyun Nam, Piotr Dollár, and Joon Hee Han. 2014. Local decorrelation for improved pedestrian detection. In *NIPS*. 424–432.

[25] Yanwei Pang, Jin Xie, Muhammad Haris Khan, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Ling Shao. 2019. Mask-Guided Attention Network for Occluded Pedestrian Detection. In *ICCV*. 4967–4975.

[26] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. (2017).

[27] Jimmy Ren, Xiaohao Chen, Jianbo Liu, Wenxiu Sun, Jiahao Pang, Qiong Yan, Yu-Wing Tai, and Li Xu. 2017. Accurate single stage detector using recurrent rolling convolution. In *CVPR*. 5420–5428.

[28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*. 91–99.

[29] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[30] Paul Viola and Michael Jones. 2001. Rapid object detection using a boosted cascade of simple features. In *CVPR*. I–I.

[31] Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian Sun, and Chunhua Shen. 2018. Repulsion loss: Detecting pedestrians in a crowd. In *CVPR*. 7774–7783.

[32] Jialian Wu, Chunluan Zhou, Ming Yang, Qian Zhang, Yuan Li, and Junsong Yuan. 2020. Temporal-Context Enhanced Detection of Heavily Occluded Pedestrians. In *CVPR*. 13430–13439.

[33] Jiaolong Xu, Sebastian Ramos, David Vázquez, and Antonio M López. 2014. Domain adaptation of deformable part-based models. *PAMI* 36, 12 (2014), 2367–2380.

[34] Junjie Yan, Zhen Lei, Longyin Wen, and Stan Z Li. 2014. The fastest deformable part model for object detection. In *CVPR*. 2497–2504.

[35] Liliang Zhang, Liang Lin, Xiaodan Liang, and Kaiming He. 2016. Is Faster R-CNN doing well for pedestrian detection?. In *ECCV*. Springer, 443–457.

[36] Shanshan Zhang, Rodrigo Benenson, and Bernt Schiele. 2017. Citypersons: A diverse dataset for pedestrian detection. In *CVPR*. 3213–3221.

[37] Shanshan Zhang, Rodrigo Benenson, Bernt Schiele, et al. 2015. Filtered channel features for pedestrian detection.. In *CVPR*, Vol. 1. 4.

[38] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. 2018. Occlusion-aware R-CNN: detecting pedestrians in a crowd. In *ECCV*. 637–653.

[39] Chengju Zhou, Meiqing Wu, and Siew-Kei Lam. 2019. SSA-CNN: Semantic Self-Attention CNN for Pedestrian Detection. *arXiv preprint arXiv:1902.09080* (2019).