# Mutual-Supervised Feature Modulation Network for Occluded Pedestrian Detection

Ye He, Chao Zhu*, Xu-Cheng Yin
School of Computer and Communication Engineering,
University of Science and Technology Beijing
Beijing, China
Email: s20190676@xs.ustb.edu.cn, chaozhu@ustb.edu.cn, xuchengyin@ustb.edu.cn

*Abstract*—State-of-the-art pedestrian detectors have achieved significant progress on non-occluded pedestrians, yet they are still struggling under heavy occlusions. The recent occlusion handling strategy of popular two-stage approaches is to build a two-branch architecture with the help of additional visible body annotations. Nonetheless, these methods still have some weaknesses. Either the two branches are trained independently with only score-level fusion, which cannot guarantee the detectors to learn robust enough pedestrian features. Or the attention mechanisms are exploited to only emphasize on the visible body features. However, the visible body features of heavily occluded pedestrians are concentrated on a relatively small area, which will easily cause missing detections. To address the above issues, we propose in this paper a novel Mutual-Supervised Feature Modulation (MSFM) network, to better handle occluded pedestrian detection. The key MSFM module in our network calculates the similarity loss of full body boxes and visible body boxes corresponding to the same pedestrian so that the full-body detector could learn more complete and robust pedestrian features with the assist of contextual features from the occluding parts. To facilitate the MSFM module, we also propose a novel two-branch architecture, consisting of a standard full body detection branch and an extra visible body classification branch. These two branches are trained in a mutual-supervised way with full body annotations and visible body annotations, respectively. To verify the effectiveness of our proposed method, extensive experiments are conducted on two challenging pedestrian datasets: Caltech and CityPersons, and our approach achieves superior performance compared to other state-of-the-art methods on both datasets, especially in heavy occlusion cases.

## I. INTRODUCTION

Pedestrian detection is a challenging computer vision task that has been widely applied in numerous applications, such as autonomous driving, robotics, and intelligent video surveillance. State-of-the-art pedestrian detectors [1]–[6] have achieved significant progress on non-occluded pedestrians, yet they are still confronted by heavy occlusions. For example, when walking on a street, pedestrians are likely to be occluded by other pedestrians or other objects, such as cars, buildings, and bicycles. Therefore, it remains as one of the most challenging issues for a pedestrian detection approach to robustly detect partially or heavily occluded pedestrians.

Many existing approaches [1]–[4] employ a simple detection strategy that assumes entirely visible pedestrians when trained
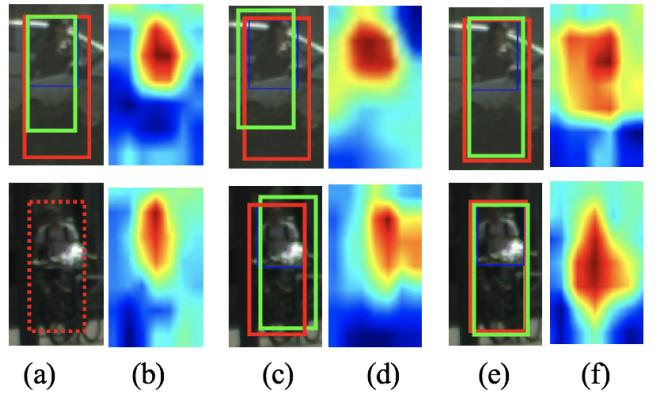
*Corresponding author



Fig. 1. Visual comparison between Bi-box, MGAN and MSFMN. (a), (c), and (e) represent the detection results of different methods. (b), (d), and (f) represent the feature visualization. Solid red boxes represent full body annotations, blue boxes are visible body annotations, green boxes denote detection results, and dashed red boxes represent the missed detections. The detected regions are cropped from the corresponding images in CityPersons val. set. Compared with Bi-box and MGAN, MSFMN displays a high response not only on the visible part but also on the occluding part.

with full body annotations. Despite achieving impressive results for non-occluded pedestrians, such a strategy is still struggling under partial or heavy occlusions since the features of the occluding part are vastly different from the visible part.

Compared with full body regions, visible parts of pedestrians usually suffer much less from occlusion, which can provide more discriminative and confident cues. Several recent approaches [7]–[9] deal with occlusions by building a two-branch architecture with extra visible-region information, available with standard pedestrian detection benchmarks, like Caltech [10] and CityPersons [11]. Nonetheless, these methods still have some weaknesses. Either the two branches are trained independently with only score-level fusion, which cannot guarantee the detectors to learn robust enough pedestrian features, such as Bi-box [7]. Or the attention mechanisms are exploited to emphasize on the visible regions while suppressing the occluded regions, like MGAN [9]. For these methods, the visible body features of the heavily occluded pedestrians are concentrated on a relatively small area, which will easily cause missing detections. Besides, some two-branch methods collect

positive training samples with full body annotations and visible body annotations simultaneously, which may sacrifice some useful visible features, such as Bi-box [7]. As illustrated in Fig. 1, (b) and (d) depict the feature maps learned by Bi-box [7] and MGAN [9], respectively. It can be observed that only the visible part has a high response while the occluded part almost has no response, and this will easily lead to inaccurate, even missing detections in heavy occlusion cases. Therefore, we think it is insufficient to only focus on the features within the visible bounding boxes in the case of heavy occlusions. The assist of the features from the occluding part is also important as a contextual cue to enhance pedestrian detectors against heavy occlusions, and this has not been studied thoroughly in previous works.

To this end, we propose in this paper a novel Mutual-Supervised Feature Modulation Network (MSFMN) aiming at enhancing feature representations of occluded pedestrians. A key part of the proposed MSFMN is the Mutual-Supervised Feature Modulation module, which aims to calculate the similarity loss between full body boxes and visible body boxes corresponding to the same pedestrian. Therefore, the full-body detector could learn more complete and robust pedestrian features in a mutual-supervised way with the assist of contextual features from the occluding parts. To obtain the visible boxes, we also construct a novel two-branch architecture consisting of a standard full body detection branch and an extra visible body classification branch. Moreover, these two branches sample their training samples supervised by full body annotations and visible body annotations, respectively (as displayed in Eq.1 and Eq.2) to obtain more focused training samples as shown in Fig. 2. Note that the proposed method can be easily applied to any existing two-stage detection framework. Finally, as shown in Fig. 1 (f), it is obvious that the feature responses of the proposed method are concentrated on a relatively wider region, including not only the visible part but also the occluding part. Therefore, the contextual information from the occluding regions can essentially enhance the discriminability of the features of occluded pedestrians.

In summary, the main contributions of this paper are three-fold:

(1) We propose a novel Mutual-Supervised Feature Modulation Network (MSFMN), to deal with the problem of occluded pedestrian detection. The key Mutual-Supervised Feature Modulation module calculates the similarity loss between full body boxes and visible body boxes to learn more robust feature representations of occluded pedestrians;

(2) The MSFMN comprises two branches: a standard full body detection branch and an extra visible body classification branch. These two branches are supervised by full body annotations and visible body annotations, respectively, to obtain more focused training samples;

(3) Extensive experiments are conducted on two standard pedestrian detection benchmarks: CityPersons [11] and Caltech [10]. Our approach sets a new state-of-the-art on both datasets, strongly validating the effectiveness of the proposed method for occluded pedestrian detection.
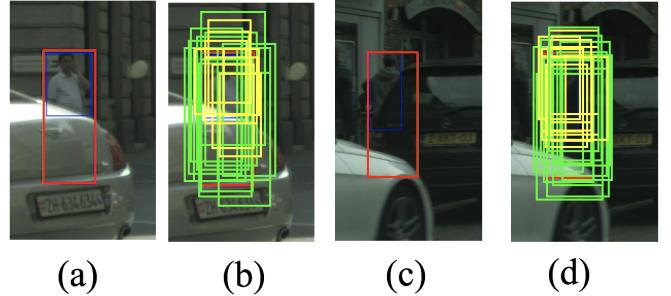


Fig. 2. Visualization of training samples obtained by our proposed sampling method. (a) and (c) are training images in CityPersons train. set. The Red boxes denote full body ground truth box and the blue boxes represent the visible body ground truth box of a pedestrian. (b) and (d) depict the positive training samples. Green boxes are positive training samples collected from full body branch and yellow boxes denote positive training samples collected from visible body branch.

## II. RELATED WORK

### A. Deep Pedestrian Detection

With the rapid development of convolutional neural networks (CNNs) [12]–[14], great progress has been made in the pedestrian detection field. Most existing CNN-based pedestrian detectors employ either one-stage or two-stage strategy as their backbone architecture. One-stage approaches [1], [2], [15], [16] where proposal generation and classification are formulated as a single-stage regression problem aim to accelerate the inference process of detectors, to meet the requirement of time efficiency in diverse real-world applications. In contrast to one-stage approaches, two-stage detectors aim to pursue the state-of-the-art performance by separate proposal generation followed by confidence computation of proposals. In recent years, two-stage pedestrian detection approaches [3]–[5], [8], [11], [17]–[19] have shown superior performance on standard pedestrian benchmarks. For example, in [17], RPN is employed to generate proposals and provide CNN features followed by a boosted decision forest. Zhang et al. [11] apply five key strategies to adapt the plain Faster R-CNN for pedestrian detection. Due to their superior performance on some pedestrian benchmarks [11], we also deploy a two-stage detection method as a backbone pipeline in this work.

### B. Occlusion Handling in Pedestrian Detection

Many efforts have been made to handle occlusions for pedestrian detection. A common strategy [6], [19]–[23] is the part-based approach where a set of part detectors are learned with each part designed to handle a specific occlusion pattern. The parts used in these approaches are usually manually designed, which may not be optimal.

Different from the above approaches, there are also some other approaches [6], [24]–[26] for occlusion handling without using parts information. In [24], an implicit shape model is adopted to generate a set of pedestrian proposals which are further refined by exploiting local and global cues. Repulsion Loss [6] and AggLoss [26] design two novel regression losses to generate more compact proposals to make them less
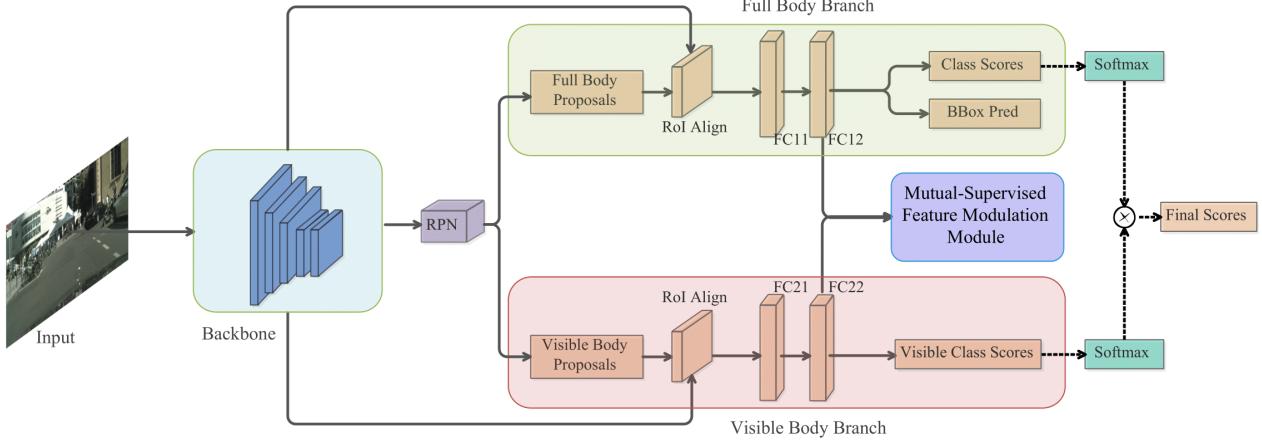
Fig. 3. Overall network architecture of our Mutual-Supervised Feature Modulation Network (MSFMN). It consists of a full body (FB) branch enclosed in the green box and a visible body (VB) branch in the red box. A novel Mutual-Supervised Feature Modulation Module is enclosed in the purple box. Two feature vectors are obtained from fully connected layer FC12 and FC22 respectively and then sent to the Mutual-Supervised Feature Modulation Module. In the architecture, the FB branch is a standard pedestrian detector branch and the VB branch is proposed to generate classification scores for visible proposals. $FC_{ij}$ denotes the j-$th$ FC layer in the i-$th$ branch. The dotted lines depict the inference process. $\otimes$ represents element-wise product operation.

sensitive to the NMS threshold. Besides, in [27], an adaptive NMS strategy is introduced that applies a dynamic suppression threshold to an instance in crowded scenes.

Contrary to the aforementioned methods, recent approaches focus on utilizing annotations of the visible body as extra supervisions together with the standard full body annotations to investigate the problem of occluded pedestrian detection. Zhang et al. [8] employ visible body information along with a pre-trained body part prediction model to learn specific occlusion patterns (full, upper-body, left-body, and right-body visible). MGAN [9], a one-way supervision network, incorporates attention mechanisms into pedestrian detection using visible region supervision to emphasize the visible regions while suppressing the occluded regions. The work of Bi-box [7] regresses full and visible body of a pedestrian at the same time. However, the two branches of Bi-box [7] are trained separately with only score-level fusion, which cannot guarantee the detectors to learn robust enough pedestrian features.

In this work, we follow the idea of utilizing extra visible annotations to tackle the problem of occluded pedestrian detection. Different to [7], our proposed method effectively integrates the two branches in the feature level to obtain more discriminative and robust features. Different to [9], our proposed method adopts a mutual-supervised way to make better use of the contextual features of occluding parts, aiming at enhancing the feature representations against heavy occlusions.

## III. PROPOSED METHOD

In this paper, we propose a novel Mutual-Supervised Feature Modulation Network (MSFMN) for occluded pedestrian detection. The overall architecture of the proposed network is described in Sec.A. To obtain the most focused positive training samples, we propose a novel proposal sampling method in Sec.B. Next, we detail the design of the novel Mutual-Supervised Feature Modulation Module in Sec.C. Finally, the total loss function of multi-task prediction for end-to-end training along with a fusion method of two branches is represented during inference in Sec.D.

### A. Overall Architecture

The overall architecture of the proposed method is illustrated in Fig. 3. Note that our proposed method can be easily applied to any existing two-stage detection frameworks. For a fair comparison, we implement it on the widely used Faster R-CNN framework [28] and adopt VGG-16 [12] as the backbone which is the most commonly used backbone in pedestrian detection networks. The architecture takes a raw image as input, first deploys a pre-trained ImageNet [29] model. Then extracted feature maps are sent to the region proposal network (RPN) to generate two different sets of candidate proposals, including full body proposals and visible body proposals. For each proposal, a fixed-sized feature representation is obtained through RoI Align [30] layer. Finally, these features go through a classification network to generate predictions. Specifically, full body (FB) branch is a standard pedestrian detection branch to generate classification scores and regressed bounding box coordinates, and for the visible body (VB) branch, we need to consider the choice of employing classification task only or both classification and regression tasks. Ref. [31] discussed that classification needs translation invariant feature whereas regression needs translation covariant feature. If we design both classification and regression tasks in the VB branch,
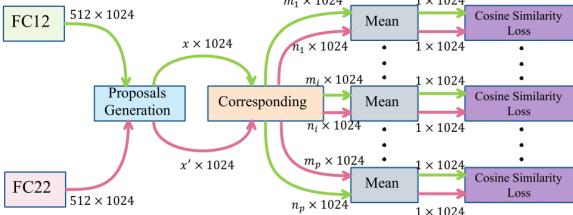
Fig. 4. The architecture of Mutual-Supervised Feature Modulation Module. It takes two 512 × 1024 feature vectors from FC12 and FC22, respectively. First, we collect $x$ and $x'$ positive training samples. Then, we find $m_i$ and $n_i$ training samples corresponding to the same pedestrian ground truth. P denotes the total number of pedestrians. Finally, we calculate cosine similarity loss of two 1 × 1024 feature vectors after $Mean$ operation.

the regression task will force the detector to gradually learn translation covariant features during training, which might potentially downgrade the performance of the classifier. Therefore, we only place a classification task in the VB branch to accurately classify visible proposals. See Sec.IV-C for more comparative results.

### B. Proposals Generation of RPN

In our architecture, the FB branch collects training samples using full body annotations and then passes them through a classification network (FC11, FC12) to generate the classification scores and the regressed bounding box coordinates. VB branch collects training samples employing visible region annotations and then feeds them into a simple classification network (FC21, FC22) to obtain the classification scores, which indicate the probability that this visible proposal contains a pedestrian. The Proposals Generation process of RPN can be expressed with the following equations:

$$P_{FB} = \{x|IoU(x, GT_{FB}) > 0.5\}\cup$$
$$\{y|IoU(y, GT_{FB}) <= 0.5\} \qquad (1)$$

$$P_{VB} = \{x'|IoU(x', GT_{VB}) > 0.5\}\cup$$
$$\{y'|IoU(y', GT_{VB}) <= 0.5\} \qquad (2)$$

where $P_{FB}$ and $P_{VB}$ represent the training samples collected for FB branch and VB branch respectively, $x$ and $x'$ denote positive training samples, and $y$ and $y'$ are negative training samples, respectively. For each branch, we sample 512 region proposals, the positive and negative samples are randomly sampled at a ratio of 1:3, following the same parameters as in [28].

### C. Mutual-Supervised Feature Modulation Module

As shown in Fig. 4, FC12 and FC22 are two fully connected layers of FB branch and VB branch, respectively. The 512 × 1024 feature vectors are passed to the Proposals Generation process stated in Sec.B to collect $x$ and $x'$ positive training samples. Then for two branches, we obtain $m_i$ and $n_i$ positive

samples corresponding to the i-$th$ pedestrian ground truth. P represents the total number of pedestrians. To measure the distance between these $m_i$ and $n_i$ samples, we compare several different methods including Manhattan distance, Euclidean distance, and Cosine Similarity. Among these measurements, Cosine Similarity achieves the best results. For simplicity, we get two 1 × 1024 feature vectors by $Mean$ operation and then calculate the cosine similarity loss of these two vectors. See Sec.IV-C for more comparative results.

The cosine similarity loss is computed as:

$$L_{MSFMM} = \frac{1}{P} \sum_{i=1}^{P}[1 - cos(\frac{1}{m_i} \sum v_i, \frac{1}{n_i} \sum v'_i)] \qquad (3)$$

where $v_i$ and $v'_i$ represent the $m_i$ × 1024 and $n_i$ × 1024 feature vectors of the i-$th$ pedestrian extracted from FB and VB branch, respectively.

This mutual-supervised loss function is intended to enhance the feature representations by incorporating the visible part features and occluding part features simultaneously.

### D. Multi-task Optimization & Inference

Here, we present the loss function for the proposed architecture MSFMN. The overall loss formulation $L$ is as follows:

$$L = L_{RPN_{cls}} + L_{RPN_{reg}} + L_{FB_{cls}}$$
$$+L_{FB_{reg}} + L_{VB_{cls}} + L_{MSFMM} \qquad (4)$$

where $L_{RPN_{cls}}$ and $L_{RPN_{reg}}$ refer to the classification and regression loss of RPN, $L_{FB_{cls}}$ and $L_{VB_{cls}}$ refer to the classification loss of FB and VB, $L_{FB_{reg}}$ is the bounding box regression loss of FB and $L_{MSFMM}$ is the loss of Mutual-Supervised Feature Modulation Module. Here, classification loss is Cross-Entropy loss and the bounding box regression loss is Smooth-L1 loss.

In the inference stage, we propose a method to fuse the information of the two branches. Since the visible box contains the most discriminative information of pedestrians, taking the classification score of the VB branch as a part of the final score of detection box will improve the accuracy of pedestrian detector. Specifically, the classification scores of the VB branch are multiplied by those of the FB branch as the scores of the final detection box. Formally, the final scores of a pedestrian are defined as:

$$Final\ Scores = Softmax(classification\ scores)\otimes$$
$$Softmax(visible\ classification\ scores) \qquad (5)$$

where $classification\ scores$ and $visible\ classification\ scores$ represent the raw scores output from FC12 and FC22, respectively, $\otimes$ denotes element-wise product operation.

## IV. EXPERIMENTAL EVALUATION

In this section, we evaluate the proposed approach on CityPersons [11] and Caltech [10], which are two challenging pedestrian detection datasets with different occlusion settings.

| Method | VB Branch | MSFMM | R | HO |
|---|---|---|---|---|
| Baseline | × | × | 11.92 | 47.88 |
| **Our MSFMN** | cls+reg | × | 11.45 | 45.35 |
| | cls | × | 10.78 | 44.81 |
| | cls | pos+neg | 10.52 | 41.35 |
| | cls | pos | **10.12** | **38.45** |

## A. Datasets and Evaluation Metrics

**Datasets:** CityPersons [11] consists of 2975 training, 500 validation, and 1525 test images. It is a challenging dataset for pedestrian detection, which exhibits large diversity. Caltech pedestrian is a popular dataset [10] containing 11 sets of videos. The first six video sets S0-S5 are for training and the remaining five video sets S6-S10 are used for testing. To increase the size of training set, we train the model on Caltech10×. Finally, the training and test sets have 42782 and 4024 images, respectively. Both datasets provide box annotations for full body and visible region.

**Evaluation Metrics:** We report the performance using log-average miss rate (MR) throughout the experiments. It is computed over the false positive per image (FPPI) range of $[10^{-2}, 10^{0}]$ [10], the lower value represents better detection performance. On Cityperons, we follow [11] and report the results across two different subsets: Reasonable (**R**), Heavy Occlusion (**HO**). For the Caltech dataset, we report results on Reasonable (**R**), Heavy Occlusion (**HO**), and the combined Reasonable + Heavy Occlusion (**R+HO**). The visibility ratio in **R** set is larger than 65%, and the visibility ratio in **HO** set ranges from 20% to 65%. Thus, the visibility ratio in **R+HO** set is larger than 20%. In all subsets, the height of pedestrians over 50 pixels is taken for evaluation, as in [8]. Notice that **HO** set is designed to evaluate the performance under severe occlusions.

## B. Implementation Details

For both datasets, the network is trained on two GPUs with a total of 2 images per mini-batch. We adopt RoI Align [30] instead of RoI Pooling for feature extraction to get more precise features. We now detail settings specific to the two datasets.

**Citypersons**. We fine-tune pre-trained ImageNet VGG model [12] on the trainset of the CityPersons. We follow the same experimental protocol as in [11] and employ two fully connected layers with 1024 instead of 4096 output dimensions. We choose SGD with momentum of 0.9 as the optimizer and set the initial learning rate as 0.0025. We train 15 epochs in total and decrease the learning rate by 0.1 at the $8$-$th$ and $11$-$th$ epochs. As in [7], ground-truth pedestrian examples which are at least 50 pixels tall and are occluded less than 70% are used for training.

**Caltech**. We start with a model pre-trained on CityPersons dataset. The initial learning rate is 0.0025 for the first 3 epochs and is reduced by 10 and 100 times for another 2 and 1 epochs.

Multi-scale training and testing are not applied to ensure fair comparisons with previous methods.

## C. Ablation Study

To sufficiently verify the effectiveness of the proposed components, we conduct detailed ablation studies on CityPersons dataset.

**Baseline Comparison.** Tab. I shows the performance of baseline and our proposed method on CityPersons validation subsets. For a fair comparison, we use the same training data, input scale (×1.3), and network backbone (VGG-16). The best results are boldfaced and shown in the final row. The baseline detector obtains a log-average miss rate of 11.92% on **R** set of CityPersons dataset, outperforming the adapted FasterRCNN baseline in CityPersons [11] by 0.89%. Thus, our baseline models are strong enough to verify the effectiveness of the proposed components. To analyze the contributions of the proposed components individually, we gradually apply the VB branch and Mutual-Supervised Feature Modulation Module (MSFMM) to the baseline model. As shown in Tab. I, our final MSFMN significantly reduces the miss rates on both **R** and **HO** subsets. Under heavy occlusions (**HO**), MSFMN achieves an absolute reduction of 9.43% in log-average miss rate compared to the baseline, demonstrates the effectiveness of MSFMN towards handling heavy occlusions.

The relevant ablation studies and analyses are presented in the following.

**Influence of VB branch.** To evaluate the efficacy of the proposed VB branch, we first add a VB branch based on the baseline model. We not only explore the effect of using classification task only in the VB branch but also explore the effect of using classification and regression tasks simultaneously. The results in Tab. I validate our analyses in Sec.III-A. Specifically, the two branches network with only a classification task in VB branch outperforms baseline model on both **R** and **HO** sets by achieving a log-average miss rate of 10.78% and 44.81%, respectively.

**Influence of different proposals for Mutual-Supervised Feature Modulation Module.** We then apply Mutual-Supervised Feature Modulation Module for two-branch detector consisting of FB branch and VB branch (with classification task only) to demonstrate its effectiveness. We not only consider applying similarity loss for positive samples but also consider applying similarity loss for both positive and negative samples. However, negative proposals have low similarity since they usually contain vastly different features. Therefore, it is more reasonable to calculate similarity loss for positive samples. Tab. I confirms that calculating similarity loss of positive samples achieves a log-average miss rate of 10.12% and 38.45% on **R** and **HO**, respectively.

**Influence of different similarity loss functions.** Next, we investigate the effect of different similarity loss functions, including Manhattan distance, Euclidean distance, and Cosine Similarity. The Manhattan distance is the distance between two points measured along axes at right angles, and Euclidean distance is the "ordinary" straight-line distance between two

TABLE II
COMPARISON (IN LOG-AVERAGE MISS RATES) OF DIFFERENT SIMILARITY METHOD.

| Method | R | HO |
|---|---|---|
| Manhattan distance | 11.25 | 46.79 |
| Euclidean distance | 10.62 | 40.61 |
| Cosine Similarity | **10.12** | **38.45** |

TABLE III
COMPARISON (IN LOG-AVERAGE MISS RATES) WITH STATE-OF-THE-ART ON THE CITYPERSONS VAL. SET.

| Method | Backbone | R | HO |
|---|---|---|---|
| Adaptive Faster RCNN [11] | VGG-16 | 12.81 | - |
| Rep.Loss [6] | ResNet-50 | 11.60 | 55.3 |
| Bi-box [7] | VGG-16 | 11.24 | 44.15 |
| FRCN+A+DT [32] | VGG-16 | 11.10 | 44.30 |
| OR-CNN [26] | VGG-16 | 11.00 | 51.30 |
| Adaptive NMS [27] | VGG-16 | 10.80 | 54.00 |
| MGAN [9] | VGG-16 | 10.50 | 39.40 |
| **Our MSFMN** | VGG-16 | **10.12** | **38.45** |

points in Euclidean space. And Cosine Similarity is to measure the difference between two individuals by cosine value of the angle between two vectors in vector space. As shown in Tab. II, Cosine Similarity method outperforms both Manhattan distance and Euclidean distance methods on both **R** and **HO** sets by achieving a log-average miss rate of 10.12% and 38.45%, respectively.

*D. State-of-the-art Comparison on CityPersons*

We compare our method with other recent state-of-the-art methods including Adapted FasterRCNN [11], Rep. Loss [6], OR-CNN [26], Bi-box [7], Adaptive NMS [27], FRCN+A+DT [32] and MGAN [9] on CityPersons dataset. As shown in Tab. III , we report the performance of MSFMN and other methods on the validation set using the same ground-truth pedestrian examples and input scale during training. The proposed MSFMN outperforms all the other methods on **R** and **HO** subsets. Notably, on **HO**, our method reduces the MR of state-of-the-art result from 39.40% to 38.45%, demonstrating the superiority of the proposed method in heavy occlusion cases. Fig. 5 displays example detections from Bi-box [7], MGAN [9], and the proposed MSFMN on CityPersons val. set. The occlusion degrees of examples vary widely from partial to heavy occlusion. Our MSFMN accurately detects pedestrians with varying levels of occlusions.

*E. State-of-the-art Comparison on Caltech*

Finally, we evaluate the MSFMN on Caltech [10] and compare it with state-of-the-art approaches. Tab. IV shows the comparison on Caltech test set under three occlusion subsets: **R**, **HO**, and **R + HO**. $*^O$ means the result is under the standard (old) test annotations, and $*^N$ means the result is under the new annotations provided by [33]. Compared to existing methods, the MSFMN achieves superior detection performance on all these subsets with log-average miss rate of 6.45%, 38.01%, 13.40%, and 2.80%, respectively. Fig. 6

TABLE IV
COMPARISON (IN LOG-AVERAGE MISS RATES) WITH STATE-OF-THE-ART COMPARISON ON CALTECH TEST SET.

| Method | $R^O$ | $HO^O$ | $R+HO^O$ | $R^N$ |
|---|---|---|---|---|
| DeepParts [21] | 11.89 | 60.42 | 22.79 | 12.90 |
| MS-CNN [3] | 9.95 | 59.94 | 21.53 | 8.08 |
| ATT-part [8] | 10.33 | 45.18 | 18.21 | 8.11 |
| SDS-RCNN [5] | 7.36 | 58.55 | 19.72 | 6.44 |
| OR-CNN [26] | - | - | - | 4.10 |
| Rep.Loss [6] | - | - | - | 4.00 |
| Bi-box [7] | 7.61 | 44.40 | 16.06 | - |
| MGAN [9] | 6.83 | 38.16 | 13.84 | - |
| **Our MSFMN** | **6.45** | **38.01** | **13.40** | **2.80** |

depicts the detection examples of the proposed MSFMN and Bi-box [7] and MGAN [9]. Our proposed MSFMN provides more accurate detections under different occlusion scenarios compared to the other two approaches.

## V. CONCLUSION

We propose a novel Mutual-Supervised Feature Modulation Network (MSFMN) for occluded pedestrian detection. A new Mutual-Supervised Feature Modulation Module is designed to calculate the similarity loss of full body boxes and visible body boxes corresponding to the same pedestrian aiming at enhancing feature representations of heavily occluded pedestrians. Our MSFMN consists of two branches, a standard full body detection branch and an extra visible body classification branch. Moreover, we utilize the full body annotations and visible annotations to supervise the two branches respectively. The effectiveness of the proposed MSFMN is validated on the CityPersons and Caltech datasets. Experimental results demonstrate that the proposed MSFMN outperforms other state-of-the-art approaches, validating the effectiveness of the proposed method.

**Future work**: the final score generation in the inference stage is relatively simple, therefore we plan to explore more effective ways to combine the information of two branches in the future.

## REFERENCES

[1] Jimmy Ren, Xiaohao Chen, Jianbo Liu, Wenxiu Sun, Jiahao Pang, Qiong Yan, Yu-Wing Tai, and Li Xu. Accurate single stage detector using recurrent rolling convolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5420–5428, 2017.
[2] Wei Liu, Shengcai Liao, Weidong Hu, Xuezhi Liang, and Xiao Chen. Learning efficient single-stage pedestrian detectors by asymptotic localization fitting. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 618–634, 2018.
[3] Zhaowei Cai, Quanfu Fan, Rogerio S Feris, and Nuno Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In *European conference on computer vision*, pages 354–370. Springer, 2016.
[4] Jiayuan Mao, Tete Xiao, Yuning Jiang, and Zhimin Cao. What can help pedestrian detection? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3127–3136, 2017.

Fig. 5. Qualitative detection examples using (a) Bi-box [7], (b) the state-of-the-art MGAN [9] and (c) MSFMN on CityPersons val. images. Solid red boxes denote the ground-truth, dashed red boxes represent the missed detections and detector predictions are indicated by green boxes. The detected regions are cropped from the corresponding images for improved visualization. Note that all detection results are obtained using the same false positive per image (FPPI) criterion. Our MSFMN accurately detects pedestrians with varying levels of occlusions.

[5] Garrick Brazil, Xi Yin, and Xiaoming Liu. Illuminating pedestrians via simultaneous detection & segmentation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4950–4959, 2017.

[6] Xinlong Wang, Tete Xiao, Yuning Jiang, Shuai Shao, Jian Sun, and Chunhua Shen. Repulsion loss: Detecting pedestrians in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7774–7783, 2018.

[7] Chunluan Zhou and Junsong Yuan. Bi-box regression for pedestrian detection and occlusion estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 135–151, 2018.

[8] Shanshan Zhang, Jian Yang, and Bernt Schiele. Occluded pedestrian detection through guided attention in cnns. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6995–7003, 2018.

[9] Yanwei Pang, Jin Xie, Muhammad Haris Khan, Rao Muhammad Anwer, Fahad Shahbaz Khan, and Ling Shao. Mask-guided attention network for occluded pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4967–4975, 2019.

[10] Piotr Dollar, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE transactions on pattern analysis and machine intelligence*, 34(4):743–761, 2011.

[11] Shanshan Zhang, Rodrigo Benenson, and Bernt Schiele. Citypersons: A diverse dataset for pedestrian detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3213–3221, 2017.

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

[13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[14] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7132–7141, 2018.

[15] Chunze Lin, Jiwen Lu, Gang Wang, and Jie Zhou. Graininess-aware deep feature learning for pedestrian detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 732–747, 2018.

[16] Junhyug Noh, Soochan Lee, Beomsu Kim, and Gunhee Kim. Improving occlusion and hard negative handling for single-stage pedestrian detectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 966–974, 2018.

[17] Liliang Zhang, Liang Lin, Xiaodan Liang, and Kaiming He. Is faster r-cnn doing well for pedestrian detection? In *European conference on computer vision*, pages 443–457. Springer, 2016.

[18] Tao Song, Leiyu Sun, Di Xie, Haiming Sun, and Shiliang Pu. Small-scale pedestrian detection based on topological line localization and temporal feature aggregation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 536–551, 2018.

[19] Garrick Brazil and Xiaoming Liu. Pedestrian detection with autoregressive network phases. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7231–7240, 2019.

Fig. 6. Qualitative detection comparison of (a) Bi-box [7], (b) MGAN [9] and (c) MSFMN under different occlusions on caltech test images. All detection results are obtained using the same false positive per image criterion. The solid red boxes denote the ground-truth, dashed red boxes represent the missed detections and the green boxes present detection results. For better visualization, the detected regions are cropped from the corresponding images.

[20] Chunluan Zhou and Junsong Yuan. Non-rectangular part discovery for object detection. In *BMVC*, 2014.

[21] Yonglong Tian, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning strong parts for pedestrian detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1904–1912, 2015.

[22] Wanli Ouyang and Xiaogang Wang. Joint deep learning for pedestrian detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2056–2063, 2013.

[23] Chunluan Zhou and Junsong Yuan. Multi-label learning of part detectors for heavily occluded pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3486–3495, 2017.

[24] Bastian Leibe, Edgar Seemann, and Bernt Schiele. Pedestrian detection in crowded scenes. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 878–885. IEEE, 2005.

[25] Xiaoyu Wang, Tony X Han, and Shuicheng Yan. An hog-lbp human detector with partial occlusion handling. In *2009 IEEE 12th international conference on computer vision*, pages 32–39. IEEE, 2009.

[26] Shifeng Zhang, Longyin Wen, Xiao Bian, Zhen Lei, and Stan Z Li. Occlusion-aware r-cnn: detecting pedestrians in a crowd. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 637–653, 2018.

[27] Songtao Liu, Di Huang, and Yunhong Wang. Adaptive nms: Refining pedestrian detection in a crowd. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6459–6468, 2019.

[28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015.

[29] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.

[30] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.

[31] Bowen Cheng, Yunchao Wei, Honghui Shi, Rogerio Feris, Jinjun Xiong, and Thomas Huang. Revisiting rcnn: On awakening the classification power of faster rcnn. In *Proceedings of the European conference on computer vision (ECCV)*, pages 453–468, 2018.

[32] Chunluan Zhou, Ming Yang, and Junsong Yuan. Discriminative feature transformation for occluded pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9557–9566, 2019.

[33] Shanshan Zhang, Rodrigo Benenson, Mohamed Omran, Jan Hosang, and Bernt Schiele. How far are we from solving pedestrian detection. pages 1259–1267, 2016.