

---

# Pedestrian Detection in Aerial Image Based on Convolutional Neural Network with Attention Mechanism and Multi-scale Prediction

---

Jiaxi Yang , [Qian Zhang](#) , Yuhang Chen , [Sitao Luan](#) \*

Posted Date: 23 January 2024

doi: 10.20944/preprints202401.1672.v1

Keywords: pedestrian detection; aerial image; attention mechanism; multi-scale prediction; convolutional neural network; new benchmark dataset



Preprints.org is a free multidiscipline platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This is an open access article distributed under the Creative Commons Attribution License which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Article

# Pedestrian Detection in Aerial Image Based on Convolutional Neural Network with Attention Mechanism and Multiscale Prediction

Jiayi Yang<sup>1</sup>, Qian Zhang<sup>1</sup>, Yuhang Chen<sup>1</sup> and Sitao Luan<sup>2,\*</sup>

<sup>1</sup> Department of Electrical and Computer, Engineering, Concordia University, Montreal, Canada; ya\_jiayi@live.concordia.ca (J.Y.); qian.zhang.20231@mail.concordia.ca (Q.Z.); yuhang.chen@mail.concordia.ca (Y.C.)

<sup>2</sup> School of Computer Science, McGill University, Mila, Montreal, Canada

\* Correspondence: sitao.luan@mail.mcgill.ca

**Abstract:** Pedestrian object detection plays a significant role in intelligent systems such as intelligent traffic and monitoring. Traditional machine learning methods on pedestrian detection have shown various drawbacks, *e.g.*, low accuracy, slow speed, *etc.* The Convolutional Neural Network (CNN) based object detection algorithms have demonstrated remarkable advantages in the field of pedestrian detection. However, the mainstream CNNs still face the problems of slow speed and low detection accuracy, especially on small and occluded targets from aerial perspective. In this paper, we propose Multi-Scale Attention YOLO (MSA-YOLO) detection algorithm to address the above issues. MSA-YOLO includes a Squeeze, Excitation and Cross Stage Partial (SECS) channel attention module for CNNs to extract richer pedestrian features with a small number of extra parameters. It also contains a multi-scale prediction module to capture the information among different pedestrian scales, which can recognize the small objects with higher accuracy and significantly reduce the missed detection. To sufficiently evaluate our proposed model, we manually collect and annotate a new benchmark dataset, Luoyang Pedestrian Dataset (The dataset can be downloaded through this anonymous link: [https://drive.google.com/drive/folders/13po1FX7Qk5RDgb60-dzOi74cq\\_kqDJtM?usp=sharing](https://drive.google.com/drive/folders/13po1FX7Qk5RDgb60-dzOi74cq_kqDJtM?usp=sharing)), which has much more sample annotations, features, scenes and image view angles than the existing benchmark datasets. In addition, the images in our dataset have higher resolution than most of benchmark pedestrian detection datasets, which can provide more detailed features of pedestrians and thus improve the model performance. When tested on Luoyang Pedestrian Dataset, our proposed MSA-YOLO algorithm significantly outperform the most commonly used baseline models with almost the same model size. This shows the efficiency of our proposed model. (The code and new dataset will be released to the public later.)

**Keywords:** pedestrian detection; aerial image; attention mechanism; multi-scale prediction; convolutional neural network; new benchmark dataset

## 1. Introduction

Pedestrian detection from an aerial perspective has abundant application scenarios [1]. For example, in the traffic field, the detection of the pedestrian can identify the residents who violate traffic regulations and enhance traffic safety issues[2]. In disaster relief missions, the use of pedestrian detection technology from an aerial viewpoint can assist rescue teams in quickly locating people who are trapped or in need of assistance [3].

Nowadays, there are two mainstream methods to solve detection tasks, one is traditional machine learning [4] and the other is deep learning [5]. The former approach consists of three phases: 1) Determine the position and range of the objects in the image; 2) Feature extractors such as HOG (Histogram of oriented gradients) are used to extract features [5]; 3) Support Vector Machine (SVM) [6] is used to classify objects according to the extracted features. This algorithm is based on appearance

features, which uses the contour information of pedestrians for classification and recognition. Since pedestrian images have different scales and spatial randomness, such detection method has low accuracy and efficiency.

Since 2006, deep learning [7] has revolutionized lots of areas [8–22], including object detection [23]. Two different methods, one-stage and two-stage algorithms, are used in deep learning for object detection. The one-stage approach utilizes features extracted by convolutional neural networks for classification and bounding box regression, and it is relatively fast in detection, *e.g.*, SSD [24], YOLO [23], RetinaNet [25], *etc.* The two-stage method, which has higher detection accuracy, takes much higher computational cost [26] than the one-stage method. It begins by using the Region Proposal Network (RPN) [27] to extract the objective region and a Convolutional Neural Network (CNN) is used to categorize and identify the candidate region. R-CNN [28] and Faster R-CNN [27] belong to this category. The two-stage methods have good robustness and higher detection accuracy, but the model size and inference time will far exceed those of the single-stage YOLO algorithm.

Nowadays, deep learning based pedestrian detection algorithms still face the low detection accuracy challenge [29] and there are two main reasons for it. The first reason is that the size of the detected pedestrians varies widely in the image. In particular, the size of the pedestrians in some regions of the image will be relatively tiny. This makes the model vulnerable to the detection of small targets. The second reason is that the objects in the pictures are often accompanied by messy background information, *e.g.*, being covered by buildings, trees and other pedestrians, which leads to lots of missed detection. Besides, we also lack high-quality benchmark dataset with sufficient aerial perspective images to train the detection model.

To address the above problems, in this paper, we propose MSA-YOLO deep learning algorithm. It has the Squeeze, Excitation and Cross Stage Partial (SECS) channel attention mechanism module, which can concentrate more on the layer of features that provide the most information and exclude the information from the less significant aspects, so that the accuracy of the detected objects can be improved. In addition, MSA-YOLO includes a multi-scale prediction module to increase the capacity on the recognition of relatively small objects in the images and decreases the rate of missed detection. These two proposed methods effectively address the problems caused by the inconspicuous pedestrian features in the images and the small size of pedestrians. To sufficiently train and evaluate our model, we collect and annotate a new benchmark dataset, Luoyang Pedestrian Dataset. After comparison with baseline models on Luoyang Pedestrian Dataset, we found that MSA-YOLO significantly outperform the baselines 2.3% without adding much computational cost.

In summary, the main contributions of this paper are as follows:

- We proposed the Squeeze, Excitation and Cross Stage Partial (SECS) channel attention module, which can extract the feature more accurately and effectively.
- Then, we propose a multi-scale prediction module, which can capture multi-scale information for small and occluded pedestrians.
- To assess the pedestrian detection models, we created a new dataset, the Luoyang Pedestrian Dataset, which contains 1200 aerial images with approximately 22800 labeled samples. The advantages of our proposed dataset are the richness of image samples, high image resolution, complexity of the scene. And compared to the currently existing pedestrian detection datasets, the camera angle we used is main from aerial view, which is unique and can fill the gap in the current pedestrian detection datasets.

The paper will be organized as follows: In Section 2, we introduce current research on deep learning approaches to pedestrian detection, covering innovative strategies for facing occlusion detection and attention mechanisms; in Section 3, we introduce the principle and structure of YOLOv5 and SENet in detail, which are important skeletons for our proposed model; in Section 4, we propose the SECS Attention Module and Multi-scale Prediction Module and name our proposed algorithm as MSA-YOLO in order to enhance the feature extraction capability for small and occluded pedestrians; in Section 5, we introduce the new benchmark datasets and evaluate the capacity of MSA-YOLO.

Besides, we conduct ablation study to verify the effectiveness of our proposed method; in Section 6, we summarize the merits of MSA-YOLO and Luoyang Pedestrian Dataset.

## 2. Related Work

### 2.1. Deep Learning-based Pedestrian Detection

The authors in [30] propose an approach for detecting occluded pedestrians. It enhances visible pedestrian areas while suppressing occluded pedestrians by modifying full-body characteristics. Additionally, they describe the occlusion-sensitive hard sample mining strategy, which prioritizes detection failures in highly obstructed pedestrians by mining hard samples based on the degree of occlusion. For enhancing pedestrian identification, Hsu *et al.* [31] suggest a brand-new stationary wavelet dilation residual super-resolution (SWDR-SR) network. In order to better maintain boundary features and enhance pedestrian recognition, they also suggest a novel low-to-high frequency connection technique (L2HFC). SWDR-SR performs better in identifying small-sized pedestrian pictures compared to baseline methods. In [32], the authors describe a novel Pose-Embedding Network for pedestrian identification that combines the Pedestrian Recognition Network (PRN) and the Region Proposal Network (RPN). The functions of these two networks are to produce candidate regions, raise confidence levels, and get rid of false positives. The effectiveness of their proposed method in comparison to the state-of-the-art (SOTA) method was demonstrated using the Caltech, CityPersons, and COPpersons datasets. In [33], the authors propose a new deep small-scale sense network for pedestrian detection, which can generate the proposed areas to detect small-scale pedestrians effectively. Additionally, they add a brand-new cross-entropy loss function to boost the loss contribution of minute pedestrians, which are challenging to detect. Their method shows outstanding detection performance on both the VIP pedestrian and the Caltech pedestrian datasets. In [34], the authors propose a new multi-scale network to detect pedestrians with the most suitable feature maps at a specific scale by matching their perceptual fields to the object size and introducing an adversarial hidden network to enhance the robustness. With a detection speed that is twice as quick as the original network, their technique reaches cutting-edge performance. Luo *et al.* [35] propose Sequential Attention-based Distinct Part Modeling to get higher classification and regression accuracies. And their proposed method improves the mean average precision against baseline models by a large margin when evaluated on Caltech and Citypersons datasets. Hsu *et al.* [36] propose the Ratio and Scale Aware YOLO (RSA-YOLO) strategy to address the issue of low detection accuracy due to the high small pedestrian ratio. In addition, they use intelligent segmentation to split the image into two local images to solve the problem of large differences in aspect ratio<sup>1</sup>. In the results, the proposed method achieves superior performance on ETH and VOC 2012 comp4 datasets compared with baseline models.

### 2.2. Attention Mechanism for Pedestrian Detection

Du *et al.* [37] propose a synthetic aperture radar (SAR) object detection algorithm, and they enhance the network by including a multi-scale feature attention module (MFAM). By applying channel and spatial attention processes to the multi-scale feature maps, the MFAM can highlight the crucial information and decrease the interference brought on by clutter. The efficacy of the proposed method has been validated by significant experimental results based on the measured SAR dataset. In [38], the authors suggest a gaze tracking technique that incorporates the local and global binocular spatial attention mechanisms (LBSAM and GBSAM, respectively) into a network model. The GazeCapture dataset is used to validate the proposed strategy, and the results show that it performs significantly better compared to existing methods. Hu *et al.* [39] describe a hybrid attention method for

---

<sup>1</sup> Aspect ratio is two numbers separated by a colon, like 16:9 or 4:3, where the first number represents the width and the second number represents the height.

lung cancer picture segmentation that combines a spatial attention mechanism and a channel attention mechanism. The hybrid attention module is applied in DenseNet convolutional neural network [40], and their proposed method improves 24.61% compared to the baseline method in lung tumor medical images. In [41], the authors propose a residual channel attention module to suppress thin clouds in images and enhance ground scene details. The proposed method shows superiority against SOTA method in reconstructing rich ground scene details when tested on real and synthetic multi-cloud images. In [42], the authors introduce a new spatial pyramidal attention network (SPANet) that uses structural information and channel relations to better represent features. Experiments show that their proposed attention module has less parameters and is 2.3% higher in mAP compared with the baseline method. In [43], the authors introduce the Self-Attention Module (SAM) as part of the architecture of YOLOv3. When evaluated on on the BDD100K and KITTI datasets, their proposed method shows approximately 2.6% mAP improvement compared to the original yolov3 network.

### 3. Preliminaries

In this section, we will introduce YOLOv5 and Squeeze and Excitation Network, which are the two backbones of our proposed methods.

#### 3.1. YOLOv5

One of the most important skeletons of our proposed algorithm is YOLOv5<sup>2</sup>, which is one of the most commonly used neural network frameworks for object detection tasks. The YOLOv5 has the merits of fast detection speed and auto anchor, and it is favored in practical applications. As shown in Figure 1, YOLOv5 consists of backbone, neck and head parts. These three parts play different roles separately: backbone is to extract image features, neck is to mix and combine features, and head is to predict the results. The images will be sent to the backbone network at the initial step. If the images are not square, the border of the pictures will be filled with blank, and the size of the images will be resized to  $640 \times 640$  pixels. Then, we can obtain a feature layer after each Cross Stage Partial(CSP) module [44] that can be enhanced to learn more features, for a total of four feature layers with sizes of  $160 \times 160$ ,  $80 \times 80$ ,  $40 \times 40$  and  $20 \times 20$ . Then, after processing by the Spatial Pyramid Pooling-Fast(SPPF) module<sup>3</sup>, feature map fusion of partial features and global features is achieved.

In the next stage, the effective feature maps output from the backbone part is delivered to the Neck part of the network from CSP2, CSP3, and SPPF, respectively. The neck part is composed of Feature Pyramid Networks (FPN) [45] and Path Aggregation Network (PAN) structure[46]. The combination of these two structures can fuse the feature layers of different shapes to extract better features. Eventually, the three feature layers which are acquired in the neck part are fed into the Head part and the final results are output using the CIoU loss function [47] and the Non-Maximal Suppression (NMS) algorithm<sup>4</sup>. The YOLOv5 has three detection heads: if we input  $640 \times 640$  size images, we can get  $80 \times 80$  feature maps for detecting  $8 \times 8$  size objects,  $40 \times 40$  feature maps for detecting  $16 \times 16$  size objects and  $20 \times 20$  feature maps for  $32 \times 32$  size objects.

---

<sup>2</sup> YOLOv5 was released by Ultralytics LLC in 2020. However, there is no publication for it.

<sup>3</sup> The SPPF layer is an improved pooling layer that aggregates data at multiple scales while preserving spatial information, improving the model's capacity to distinguish targets of various sizes while optimizing computing efficiency.

<sup>4</sup> Non-Maximal Suppression is used to increase detection accuracy by reducing overlapping candidate frames and preserving only the best candidate frames that are most likely to contain the target.

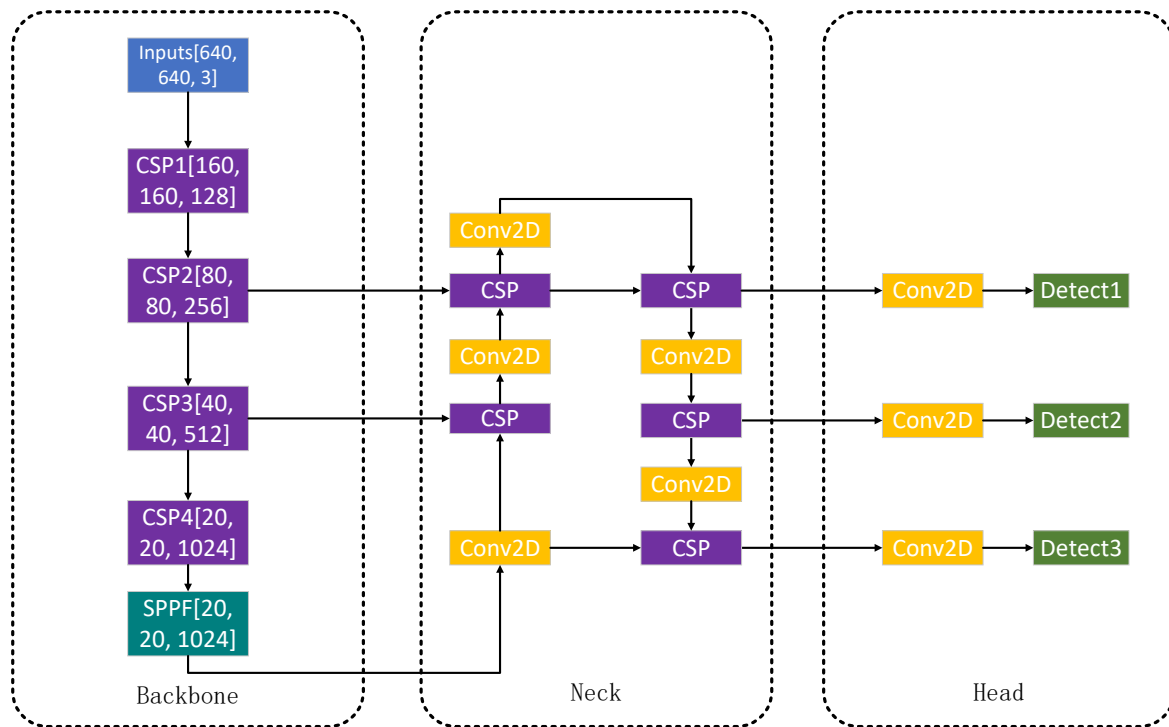


Figure 1. Structure of YOLOv5.

### 3.2. Squeeze and Excitation Network

Squeeze and Excitation Network (SENet) is a convolutional neural network which incorporates the squeeze and excitation block, *i.e.*, an attention module. The attention mechanism module only add a tiny number of extra parameters, which has very little impact on training speed and enables the network to improve model accuracy by focusing more on features that are more important to the task at hand.

As the architecture shown in Figure 2, the input feature map  $X \in \mathbb{R}^{W' \times H' \times C'}$  is fed into the header in the attention mechanism module, where  $W'$ ,  $H'$  and  $C'$  stand for the feature width, height, and number of channels. Afterward, the output of feature layer  $U = \{u_1, u_2, \dots, u_C\} \in \mathbb{R}^{W \times H \times C}$  is produced by confounding the input feature map

$$u_c = V_c * X = \sum_{s=1}^{C'} V_c^s * X^s \quad (1)$$

where  $V_c$  is the learned filter kernel set and  $*$  represents the convolution operation. The network can be made more sensitive to the information aspect with the above multi-channel convolution.

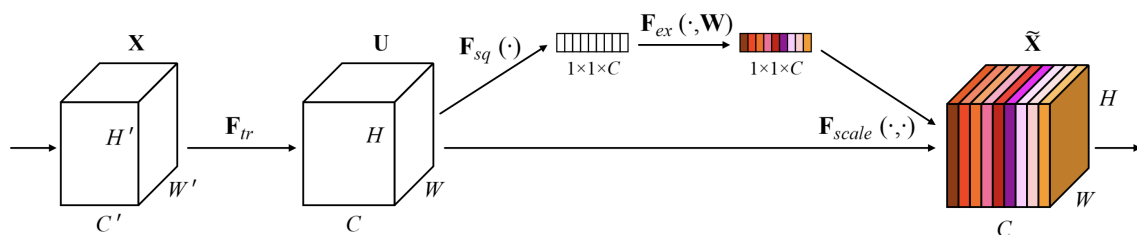


Figure 2. Squeeze and Excitation block.

In the next step, the Squeeze operation is performed on the feature map  $U$  to turn the two-dimensional feature channel into a scalar number. In other words, the feature map  $U \in \mathbb{R}^{H \times W \times C}$  is converted into a  $1 \times 1 \times C$  output, which can have a global perceptual field to some extent. The formula is shown in Equation (2), where  $F_{sq}$  indicates the global average pooling. It not only reduces the number of parameters in the module but also avoids the negative effects of too many channels on model aggregation.

$$z_c = F_{sq}(u_c) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W u_c(i, j) \quad (2)$$

After we have obtained a feature layer of size  $1 \times 1 \times C$ , the feature layers will be fed into two fully-connected layers to learn an adaptive weight for each channel and thereby decide which channels are more important to focus on.

#### 4. The Proposed Method

In this section, inspired by the YOLOv5 network architecture and the SENet attention mechanism, we propose Multi-Scale Attention YOLO (MSA-YOLO), which contains the SECSP attention module and the multiscale prediction module.

##### 4.1. SECSP Module

The neck of the YOLOv5 is made of FPN and PAN structures, however, FPN and PAN have limited feature extraction capabilities in complex scenes such as pedestrian targets that are occluded by obstacles. In order to enhance the feature extraction capability of FPN and PAN in complex environments and focus more on the essential features of the pedestrian objects, in this section, we introduced Squeeze, Excitation and Cross Stage Partial (SECSP) channel attention module.

The Squeeze and Excitation Network (SENet) is added to the PAN structure of the network, following by the CSP layer, so that the Occluded and small-sized pedestrian features in the image can be captured more effectively. As shown in Figure 3, in SECSP module, the input feature maps first flow through a CBL layer that contains convolutional operations to extract spatial features, batch normalization to stabilize the learning process and accelerate convergence, and a Leaky ReLU activation function to introduce nonlinearities and facilitate gradient propagation. The feature map is then processed through an additional convolutional layer, while at the same time a portion of the original input feature map goes directly into another convolutional layer. These two feature streams are processed through the convolutional layer and then merged at the Concat layer to integrate different levels of information. Subsequently, the merged feature maps are fed into the SENet module, where the features are recalibrated through the channel attention mechanism to highlight important information and suppress interfering information. This pipeline is demonstrated in Figure 3, which effectively improves the detection accuracy and model robustness.

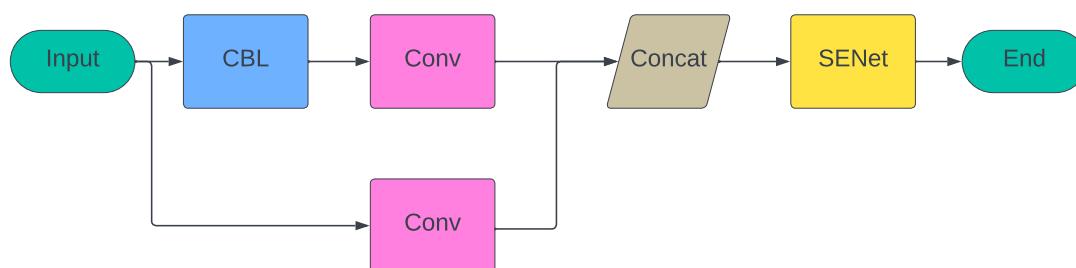


Figure 3. Structure of SECSP.

#### 4.2. Implements multi-scale prediction module for small objects

The pedestrian detection task is unique in that the size of the pedestrians in the images varies widely. The detection results of YOLOv5 convolutional neural network on the test set reveals that some pedestrians that occupy a relatively small portion of the image cannot be detected. To reduce the missed detection rate of the small-size pedestrian detection, we proposed Multi-Scale Attention YOLO (MSA-YOLO) pedestrian detection algorithm. The network structure is shown in Figure 4.

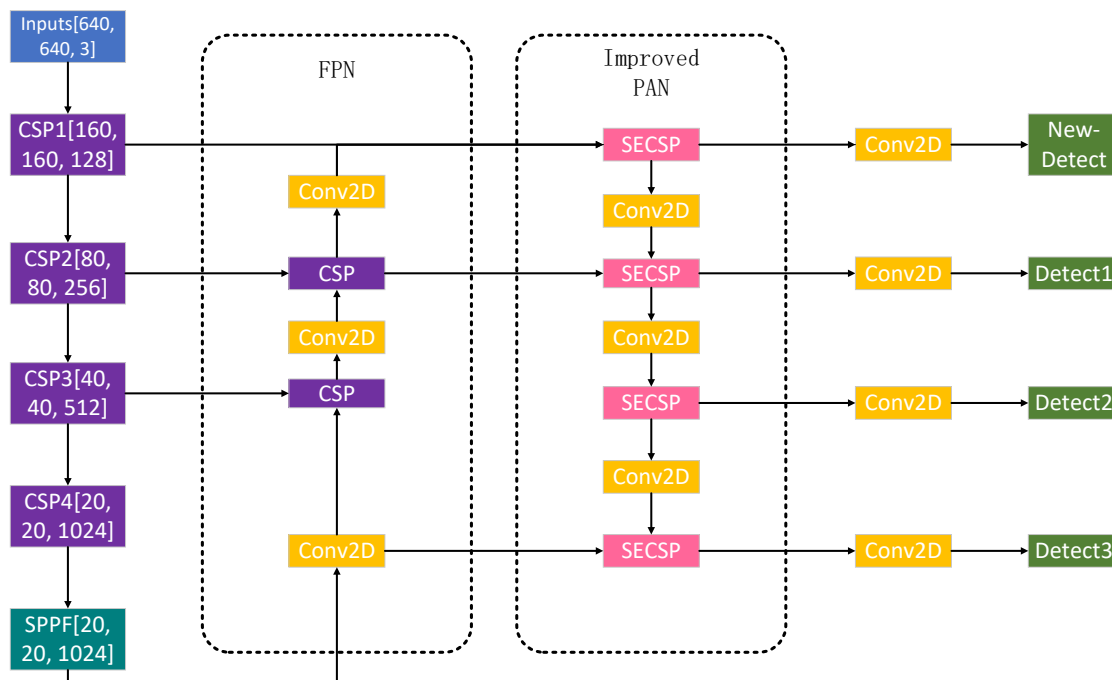


Figure 4. Structure of MSA-YOLO.

The original PAN can only output three effective feature layers that can provide three prediction scales. In addition, due to the large down-sampling multiplier and excessive perceptual field of the YOLO network, locating the feature information of small objects on deeper feature maps is quite challenging, therefore the effectiveness of small object detection is not satisfactory. To address the issues of insufficient precision and high miss detection rate brought by the scale discrepancies in the images, we introduced a multi-scale prediction module to the PAN structure. As shown in Figure 4, the new feature layer consists of the feature maps obtained from the second and third CSP layers of the backbone network at the FPN after two CSP and convolution operations, and fusion with the first CSP layer of the backbone network to obtain a feature map of  $160 \times 160$  size by using the Concat operation. The new feature layer is used to detect small objects of  $4 \times 4$ , which can effectively fuse the shallower feature maps with the deeper feature maps, thus enhancing the feature extraction capability and improving the detection accuracy of tiny targets. Although this method increases the computational cost and reduces the inference speed to some extent, the detection results are significantly improved, especially for the small targets, which are often missed by YOLOv5.

## 5. Experiments

### 5.1. Hardware, Software and Hyperparameters

In this paper, we used the Windows OS version of the Pytorch framework to build our model. The hardware device environment was an NVIDIA GeForce RTX 3090 with 24GB GPU and 64 GB

of RAM with 2933MHz. During the training of the model, we set the hyperparameters the same as YOLOv5, with a learning rate of 0.01 and a weight decay rate of 0.0005, and SGD as optimizer.

## 5.2. Dataset

We collected and created a new dataset, Luoyang Pedestrian Dataset, to evaluate the model. The samples of Luoyang Pedestrian Dataset are demonstrated in Figure 5. Our dataset elevates pedestrian detection to new heights with its aerial perspective and a high resolution of  $5472 \times 3078$  pixels, providing a level of detail unprecedented in current public datasets. Unlike the commonly found eye-level views with fewer and lower-resolution samples, our aerial dataset captures a wider scene, offering a richer array of samples for superior model training. This expansive dataset, with its abundance of annotated samples, is a robust benchmark dataset for developing advanced detection models that require detailed environmental understanding and can handle complex, real-world scenarios. Our dataset is collected by DJI Mavic Air 2S drone with three shooting angles of 35 degrees down, 45 degrees down and 55 degrees down. The dataset is composed of abundant sample types, with a total of 1200 images and about 22800 labeled samples and it contains two heights of 7.5m and 10m. This not only enriches the sample types but also enables the network to learn more pedestrian features for better application in practical scenarios.

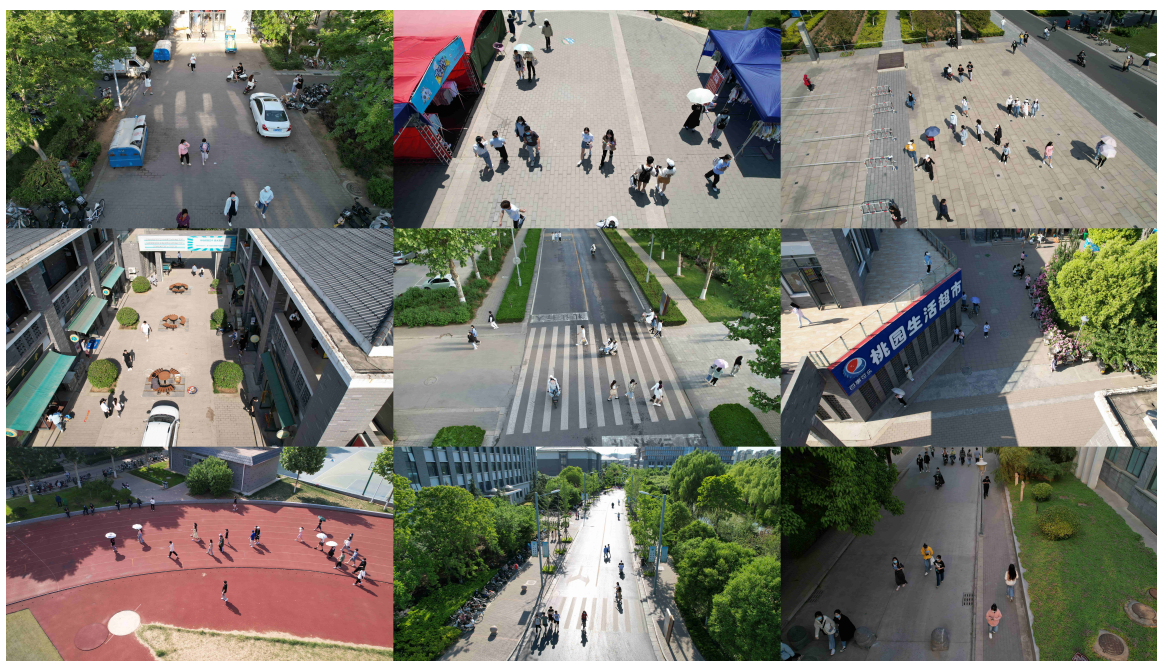


Figure 5. The samples of Luoyang Pedestrian Dataset.

## 5.3. Results and Discussion

### 5.3.1. Ablation Study

We conduct ablation experiments on Luoyang Pedestrian Dataset to verify the effectiveness of the components in our new model. The average accuracy values for each category, *i.e.*, the mAP values, are reported to assess our model. As shown in Table 1, the baseline YOLOv5 model achieves an average accuracy value of 94.7%. When the attention module is added, the accuracy increases to 95.1% ( $\uparrow$  0.4%) and when the multi-scale prediction module are included, the accuracy increases to 96.4% ( $\uparrow$  1.7%). The performance enhancement of the two partial models show the effectiveness of the attention and multi-scale prediction module. When both modules are added, the full model shows significant

improvement with an average accuracy of 97.0% ( $\uparrow$  2.3%), but the required memory is almost the same as the baseline YOLOv5. This shows the efficiency of our proposed model.

**Table 1.** Ablation Study Results.

Methods	SECSP	Multi-scale Prediction	mAP	Size
YOLOV5	×	×	94.7%	40.2M
	✓	×	95.1%	40.5M
	×	✓	96.4%	44.8M
MSA-YOLO (Ours)	✓	✓	97.0%	45.1M

### 5.3.2. Comparison with Other Baseline Models

To compare with other baseline models, we test YOLOv4 [48], Fast R-CNN [49] and Faster R-CNN [27] on Luoyang Pedestrian Dataset. <sup>5</sup> The results are reported in Table 2. It can be observed that our proposed MSA-YOLO significantly outperforms YOLOv4 with slightly larger model size, and outperforms Fast R-CNN and Faster R-CNN with significantly smaller model size. This again shows the efficiency of our proposed model.

**Table 2.** Comparison Results.

Methods	Backbone	mAP	Size
YOLOV4	Darknet53	86.7%	31.4M
YOLOV5	CSP Darknet	94.7%	40.2M
Fast R-CNN	VGG-16	95.8%	227.5M
Faster R-CNN	ResNet50	96.6%	337.1M
MSA-YOLO (Ours)	CSP Darknet	97.0%	45.1M

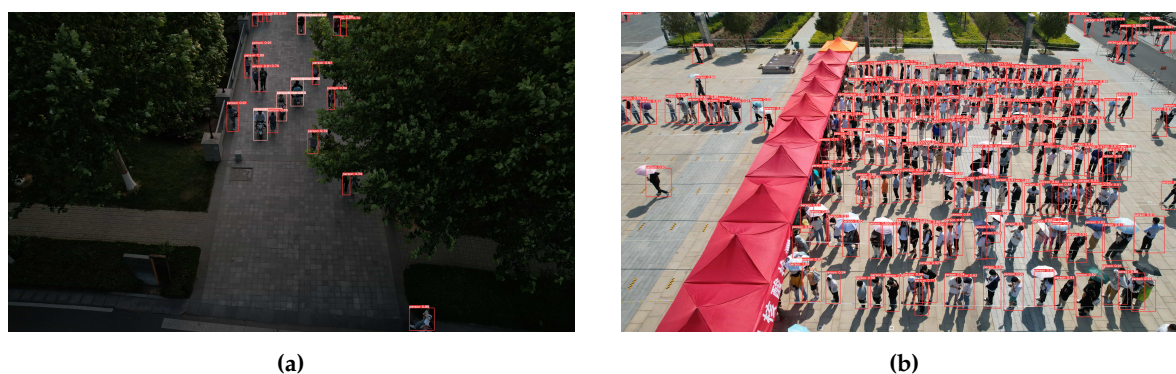
To visualize the advantages of MSA-YOLO against YOLOv5 on detecting small objects, we show the detection results of YOLOv5 and MSA-YOLO in Figure 6. The green rectangular boxes in the images represent detection errors and missed detection, and the red rectangular boxes represent correct detection results. We can find that the proposed MSA-YOLO has significantly reduced the missed detection of small objects and increased the the prediction confidence of the correct detection.

<sup>5</sup> The training of the three networks will be ended at the convergence of the model, which is the same as the original papers.



**Figure 6.** The visualization performance. (a) and (b) shows the detection results of YOLOv5, (c) and (d) shows the results of our proposed MSA-YOLO (Zoom up the figure to see the prediction confidence more clearly).

The visualization of the MSA-YOLO outputs in more complicated and difficult scenarios are shown in Figure 7. The visualization results indicate that the proposed MSA-YOLO can still detect pedestrians in the places with low light intensity, obscured by foliage and in crowded square with low detection errors and high prediction confidence.



**Figure 7.** The visualization performance. (a) shows the detection results in a place with low light intensity and obscured by foliage, and (b) shows the detection results in a crowded square (Zoom up the figure to see the prediction confidence more clearly).

## 6. Conclusions

In this paper, we propose the MSA-YOLO detection algorithm which has a stronger and lightweight attention mechanism module, SECSP, for feature extraction. In addition, a multi-scale prediction module is added to the network for the detection of small-sized objects. The combination of these two modules lead us to the proposed MSA-YOLO.

Besides, we collect and build a new dataset, Luoyang Pedestrian Dataset, which contains a great number of occluded pedestrian objects with various sizes. The ablation study, comparison with

baseline and visualization of detection results on Luoyang Pedestrian Dataset all show the efficiency of our proposed MSA-YOLO model.

## References

1. S. A. Kumar, E. Yaghoubi, A. Das, B. Harish, and H. Proenca, "The p-destre: A fully annotated dataset for pedestrian detection, tracking, and short/long-term re-identification from aerial devices," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 1696–1708, 2020.
2. A. Nguyen and B. Lee, "Enhancing traffic safety through pedestrian detection: A deep learning approach," *Journal of Transportation Safety*, vol. 29, no. 2, pp. 157–174, 2021.
3. S. Sambolek and M. Ivacic-Kos, "Automatic person detection in search and rescue operations using deep cnn detectors," *Ieee Access*, vol. 9, pp. 37905–37922, 2021.
4. M. Bilal and M. S. Hanif, "Benchmark revision for hog-svm pedestrian detector through reinvigorated training and evaluation methodologies," *IEEE transactions on intelligent transportation systems*, vol. 21, no. 3, pp. 1277–1287, 2019.
5. K. Dasgupta, A. Das, S. Das, U. Bhattacharya, and S. Yogamani, "Spatio-contextual deep network-based multimodal pedestrian detection for autonomous driving," *IEEE transactions on intelligent transportation systems*, vol. 23, no. 9, pp. 15940–15950, 2022.
6. J. Platt, "Sequential minimal optimization: A fast algorithm for training support vector machines," 1998.
7. G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
8. A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
9. A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. Ieee, 2013, pp. 6645–6649.
10. D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
11. D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.
12. M. Zhao, Z. Liu, S. Luan, S. Zhang, D. Precup, and Y. Bengio, "A consciousness-inspired planning agent for model-based reinforcement learning," *Advances in neural information processing systems*, vol. 34, pp. 1569–1581, 2021.
13. M. Zhao, S. Luan, I. Porada, X.-W. Chang, and D. Precup, "Meta-learning state-based eligibility traces for more sample-efficient policy evaluation," in *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, 2020, pp. 1647–1655.
14. T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2016.
15. S. Luan, M. Zhao, X.-W. Chang, and D. Precup, "Break the ceiling: Stronger multi-scale deep graph convolutional networks," *Advances in neural information processing systems*, vol. 32, 2019.
16. S. Luan, C. Hua, Q. Lu, J. Zhu, M. Zhao, S. Zhang, X.-W. Chang, and D. Precup, "Is heterophily a real nightmare for graph neural networks to do node classification?" *arXiv preprint arXiv:2109.05641*, 2021.
17. S. Luan, C. Hua, Q. Lu, J. Zhu, X.-W. Chang, and D. Precup, "When do we need graph neural networks for node classification?" *International Conference on Complex Networks and Their Applications*, 2023.
18. W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.
19. S. Luan, M. Zhao, C. Hua, X.-W. Chang, and D. Precup, "Complete the missing half: Augmenting aggregation filtering with diversification for graph convolutional networks," in *NeurIPS 2022 Workshop: New Frontiers in Graph Learning*, 2022.
20. S. Luan, C. Hua, Q. Lu, J. Zhu, M. Zhao, S. Zhang, X.-W. Chang, and D. Precup, "Revisiting heterophily for graph neural networks," *Advances in neural information processing systems*, vol. 35, pp. 1362–1375, 2022.

21. S. Luan, C. Hua, M. Xu, Q. Lu, J. Zhu, X.-W. Chang, J. Fu, J. Leskovec, and D. Precup, "When do graph neural networks help with node classification? investigating the impact of homophily principle on node distinguishability," *Advances in Neural Information Processing Systems*, vol. 36, 2023.
22. C. Hua, S. Luan, M. Xu, R. Ying, J. Fu, S. Ermon, and D. Precup, "Mudiff: Unified diffusion for complete molecule generation," in *The Second Learning on Graphs Conference*, 2023.
23. J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
24. W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part I 14*. Springer, 2016, pp. 21–37.
25. T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2980–2988.
26. C. Zhang and D. Kim, "Comparative analysis of one-stage and two-stage deep learning architectures for pedestrian detection," *Journal of Computer Vision*, vol. 34, no. 4, pp. 789–805, 2022.
27. S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.
28. R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.
29. S. Iftikhar, Z. Zhang, M. Asim, A. Muthanna, A. Koucheryavy, and A. A. Abd El-Latif, "Deep learning-based pedestrian detection in autonomous vehicles: Substantial issues and challenges," *Electronics*, vol. 11, no. 21, p. 3551, 2022.
30. J. Xie, Y. Pang, M. H. Khan, R. M. Anwer, F. S. Khan, and L. Shao, "Mask-guided attention network and occlusion-sensitive hard example mining for occluded pedestrian detection," *IEEE transactions on image processing*, vol. 30, pp. 3872–3884, 2020.
31. W.-Y. Hsu and P.-C. Chen, "Pedestrian detection using stationary wavelet dilated residual super-resolution," *IEEE Transactions on Instrumentation and Measurement*, vol. 71, pp. 1–11, 2022.
32. Y. Jiao, H. Yao, and C. Xu, "Pen: Pose-embedding network for pedestrian detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 3, pp. 1150–1162, 2020.
33. B. Han, Y. Wang, Z. Yang, and X. Gao, "Small-scale pedestrian detection based on deep neural network," *IEEE transactions on intelligent transportation systems*, vol. 21, no. 7, pp. 3046–3055, 2019.
34. C. Lin, J. Lu, and J. Zhou, "Multi-grained deep feature learning for robust pedestrian detection," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 12, pp. 3608–3621, 2018.
35. Y. Luo, C. Zhang, W. Lin, X. Yang, and J. Sun, "Sequential attention-based distinct part modeling for balanced pedestrian detection," *IEEE transactions on intelligent transportation systems*, vol. 23, no. 9, pp. 15 644–15 654, 2022.
36. W.-Y. Hsu and W.-Y. Lin, "Ratio-and-scale-aware yolo for pedestrian detection," *IEEE transactions on image processing*, vol. 30, pp. 934–947, 2020.
37. Y. Du, L. Du, and L. Li, "An sar target detector based on gradient harmonized mechanism and attention mechanism," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2021.
38. L. Dai, J. Liu, and Z. Ju, "Binocular feature fusion and spatial attention mechanism based gaze tracking," *IEEE Transactions on Human-Machine Systems*, vol. 52, no. 2, pp. 302–311, 2022.
39. H. Hu, Q. Li, Y. Zhao, and Y. Zhang, "Parallel deep learning algorithms with hybrid attention mechanism for image segmentation of lung tumors," *IEEE Transactions on Industrial Informatics*, vol. 17, no. 4, pp. 2880–2889, 2020.
40. G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.
41. X. Wen, Z. Pan, Y. Hu, and J. Liu, "An effective network integrating residual learning and channel attention mechanism for thin cloud removal," *IEEE Geoscience and Remote Sensing Letters*, vol. 19, pp. 1–5, 2022.
42. X. Ma, J. Guo, A. Sansom, M. McGuire, A. Kalaani, Q. Chen, S. Tang, Q. Yang, and S. Fu, "Spatial pyramid attention for deep convolutional neural networks," *IEEE Transactions on Multimedia*, vol. 23, pp. 3048–3058, 2021.

43. D. Tian, C. Lin, J. Zhou, X. Duan, Y. Cao, D. Zhao, and D. Cao, "Sa-yolov3: An efficient and accurate object detector using self-attention mechanism for autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 5, pp. 4099–4110, 2020.
44. C.-Y. Wang, H.-Y. M. Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, "Cspnet: A new backbone that can enhance learning capability of cnn," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 2020, pp. 390–391.
45. T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
46. S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8759–8768.
47. Z. Zheng, P. Wang, D. Ren, W. Liu, R. Ye, Q. Hu, and W. Zuo, "Enhancing geometric factors in model learning and inference for object detection and instance segmentation," *IEEE transactions on cybernetics*, vol. 52, no. 8, pp. 8574–8586, 2021.
48. A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao, "Yolov4: Optimal speed and accuracy of object detection," *arXiv preprint arXiv:2004.10934*, 2020.
49. R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.