

Guided Attention Network for Object Detection and Counting on Drones

Yuanqiang Cai^{1,2}, Dawei Du³, Libo Zhang^{1*}, Longyin Wen⁴, Weiqiang Wang², Yanjun Wu¹, SiweiLyu³

¹ Institute of Software Chinese Academy of Sciences, China

² University of Chinese Academy of Sciences, China

³ University at Albany, SUNY, USA

⁴ JD Digits, USA

Abstract

Object detection and counting are related but challenging problems, especially for drone based scenes with small objects and cluttered background. In this paper, we propose a new Guided Attention Network (GANet) to deal with both object detection and counting tasks based on the feature pyramid. Different from the previous methods relying on unsupervised attention modules, we fuse different scales of feature maps by using the proposed weakly-supervised Background Attention (BA) between the background and objects for more semantic feature representation. Then, the Foreground Attention (FA) module is developed to consider both global and local appearance of the object to facilitate accurate localization. Moreover, the new data argumentation strategy is designed to train a robust model in various complex scenes. Extensive experiments on three challenging benchmarks (i.e., UAVDT, CARPK and PUCPR+) show the state-of-the-art detection and counting performance of the proposed method compared with existing methods.

1. Introduction

Object detection and counting are fundamental techniques in many applications, such as scene understanding, traffic monitoring and sports video, to name a few. However, these tasks become even more challenging in drone based scenes because of various factors such as small objects, scale variation and background clutter. With the development of deep learning, much progress has been achieved recently. Specifically, deep learning based detection and counting frameworks focus on discriminative feature representation of the objects.

First of all, the feature pyramid is widely applied in deep learning because it has rich semantics at all levels, *e.g.*,

U-Net [26], TDM [27] and FPN [16]. To better exploit multi-scale feature representation, the researchers use various attention modules to fuse feature maps. In [11], the channel-wise feature responses are recalibrated adaptively by explicitly modelling interdependencies between channels. [31] propose the non-local network to capture long-range dependencies, which computes the response at a position as a weighted sum of the features at all positions. Moreover, [1] develop a lightweight global context (GC) block based on the non-local module. However, all the above methods use unsupervised attention module, but consider little about the background discriminative information in feature maps.

Based on the fused feature maps, the object is represented by proposals in anchor based methods [25, 18, 2] or key-points in anchor-free methods [14, 36, 32]. Anchor based methods exploit the global appearance information of the object, relying on pre-defined anchors. It is not flexible to design different kinds of anchors because of large scale variation in drone based scenes. Anchor-free methods employ corner points, center points or target part points to capture local object appearance without anchors. However, local appearance representation does not contain object's structure information, which is less discriminative in cluttered background, especially for small objects.

In addition, the diversity of training data is essential in deep learning. Especially in the drone based scenes, the number of difficult samples is very limited. It is difficult for traditional data argumentation such as rescale, horizontal flip, rotation and cropping to train a robust model to deal with unconstrained drone based scenarios.

To address these issues, in this paper, we propose an anchor-free Guided Attention Network (GANet). First, the background attention module can enforce different channels of feature maps to learning discriminative background information for the feature pyramid. We fuse the multi-level features with the weakly-supervision of classification between background and foreground images. Second, the foreground attention module is used to capture both global and local appearance representation of the objects by tacking

*Corresponding author: Libo Zhang(libo@iscas.ac.cn). This work is supported by the National Natural Science Foundation of China under Grant No.61807033, and the Key Research Program of Frontier Sciences, CAS, Grant No.ZDBS-LY-JSC038.

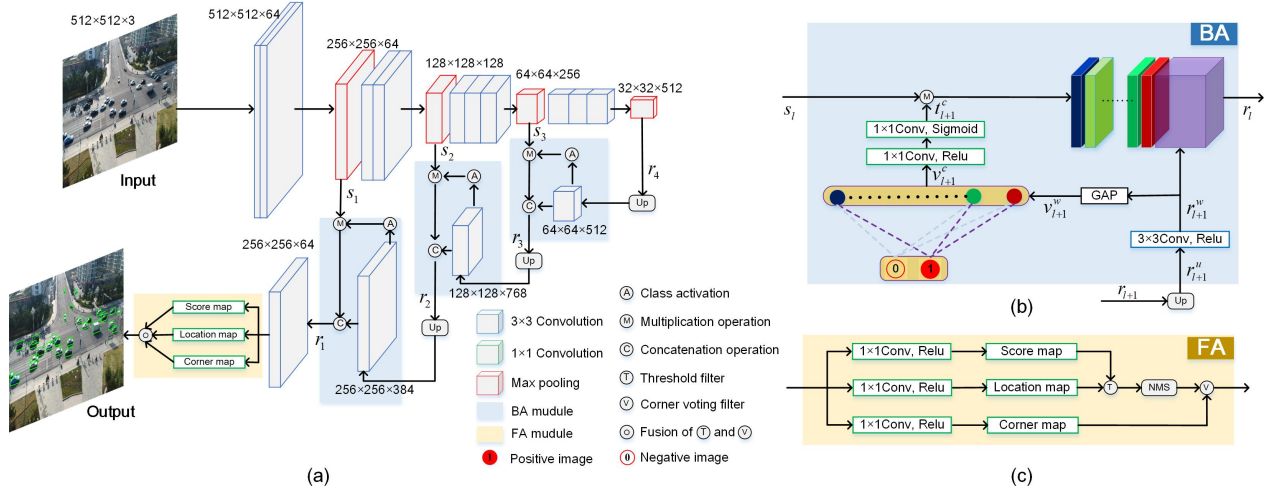


Figure 1. (a) The architecture of GANet. (b) The background attention module. (c) The foreground attention module. In (a), s_1 , s_2 , and s_3 denote *pool1*, *pool2*, and *pool3* low-level features, respectively; r_1 , r_2 , and r_3 denote the corresponding high-level features. In (b), s_l denotes the low-level features with rich texture details, r_{l+1} and r_l denote the high-level features with strong semantic information.

the merits of both anchor-based and anchor-free methods. We extract more context information in the corner regions of the object to consider local appearance information. Third, we develop a new data argumentation strategy to reduce the influence of different illumination conditions on the images for the drone based scenes, *e.g.*, sunny, night, cloudy and foggy scenes. We conduct extensive experiments on three challenging datasets (*i.e.*, UAVDT [4], CARPK [10] and PUCPR+ [10]) to show the effectiveness of the proposed method.

The main contributions of this paper are summarized as follows. (1) We present a guided attention network for object detection and counting on drones, which is formed by the foreground and background attention blocks to extract the discriminative features for accurate results. (2) A new data augmentation strategy is designed to boost up the model performance. (3) Extensive experiments on three challenging dataset, *i.e.*, UAVDT, CARPK and PUCPR+, demonstrate the favorable performance of the proposed method against the state-of-the-arts.

2. Guided Attention Network

In this section, we introduce the novel anchor-free deep learning network for object detection and counting in drone images, the Guided Attention Network (GANet), which is illustrated in Figure 1. Specifically, GANet consists of three parts, *i.e.*, the backbone, multi-scale feature fusion, and output predictor. We will first describe each part in detail, and then loss function and data argumentation strategy.

2.1. Backbone Network

Since diverse scales of objects are taken into consideration in feature representation, we choose the feature maps

from four side-outputs of the backbone network (*e.g.*, VGG-16 [28] and ResNet-50 [9]). Four side outputs correspond to *pool1*, *pool2*, *pool3*, and *pool4*, each of which is the output of four convolution blocks with different scales, respectively. The feature maps from four *pooling* layers are $\frac{1}{2}$, $\frac{1}{4}$, $\frac{1}{8}$, $\frac{1}{16}$ the size of the input image. They are marked with light blue regions in Figure 1(a). The backbone network is pre-trained by the ImageNet dataset [13].

2.2. Multi-Scale Feature Fusion

As discussed in [16], the feature pyramid has strong semantics at all scales, resulting in significant improvement as a generic feature extractor. Specifically, we fuse the side-outputs of the backbone network from top to down, *e.g.*, feature maps from *pool4* to *pool1* of VGG-16. Meanwhile, the receptive fields of the stacked feature maps can adaptively match the scale of objects. To consider background discriminative information in the feature pyramid, we introduce the Background Attention (BA) module in multi-scale feature fusion.

2.2.1 Background Attention.

As shown in Figure 1(b), the BA modules are stacked from the deepest to the shallowest convolutional layer. At the same time, the cross-entropy loss function is used to enforce different channels of feature maps focus on either foreground and background in every stage. Then, the attention module weights the pooling features with the same scale via the class-activated tensor. Finally, the weighted pooling features and the up-sampled features are concatenated and regarded as the base feature maps in the next BA.

We denote the l -th pooling features as s_l , and the input

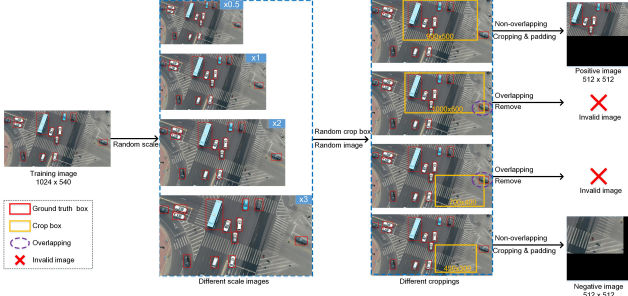


Figure 2. Generation of positive and negative samples.

and output of l -th BA as r_{l+1} and r_l . Specifically, r_{l+1} is used to learn the class-related weights for activating the class-related feature maps in s_l . For the deepest BA module, the input is regarded as the *pool4* feature maps (see r_4 in Figure 1(a)). Note that the size of output r_l in this architecture is the same as the pooling features s_l rather than the size of input r_{l+1} . Therefore, the bilinear interpolation is introduced to up-sample r_{l+1} to r_{l+1}^u . As the up-sampling operation is a linear transformation, one 3×3 convolutional layer w_l^u is used as soft-adding to improve the scale adaptability. Instead of concatenating the up-sampled r_{l+1} and the activated s_l directly, the 1×1 and 3×3 convolutional layers w_l^c is used to generate r_l . In summary, the l -th BA is formulated as

$$r_l = w_l^c \cdot (f(s_l, r_{l+1}^w) + r_{l+1}^w), \quad (1)$$

where w_l^c denotes the convolutional weights of the concatenation layer. $r_{l+1}^w = w_l^u * r_{l+1}^u$ and w_l^u are the convolutional weights of up-sampled r_{l+1}^u . w_l^c has two elements, *i.e.*, one for r_{l+1}^w and the other for $f(s_l, r_{l+1}^w)$. $f(s_l, r_{l+1}^w)$ is a class activation function with two parameters, *i.e.*, the pooling features s_l and the weighted up-sampled features r_{l+1}^w . It is defined as

$$f(s_l, r_{l+1}^w) = s_l \otimes g^c(r_{l+1}^w), \quad (2)$$

where \otimes is the multiply operation between the features s_l and the weight tensor $g^c(r_{l+1}^w)$. $g^c(r_{l+1}^w)$ is obtained by three steps. First, r_{l+1}^w is compressed into a one-dimensional vector v_{l+1}^w by the Global Average Pooling (GAP) [35]. Second, v_{l+1}^w is activated and converted to the vector with class-related information v_{l+1}^c via determining whether the input image contains the objects. Third, v_{l+1}^c is transformed into a weight tensor with class-related information $t_{l+1}^c = g^c(r_{l+1}^w)$ via two 1×1 convolutional layers.

2.2.2 Positive and Negative Image Generation.

To learn class-related feature maps, we use both the images with and without objects in the training stage. We denote them as positive and negative images respectively. Specifically, we use positive images with objects to activate the

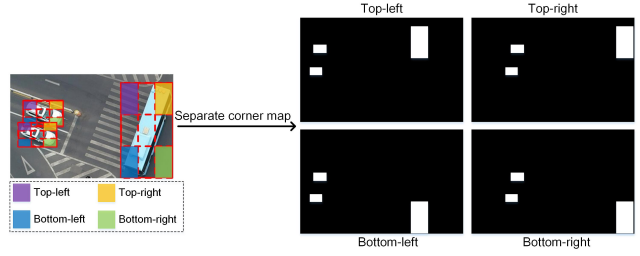


Figure 3. Illustration of four corner maps for foreground attention.

channels of feature maps to represent the pixels of object region, and negative images without overlapping of objects to activate the channels of feature maps to describe the background region. As shown in Figure 2, we generate positive and negative images with the size of 512×512 by randomly cropping and padding the rescaled training images (from 0.5x to 3x scale).

2.3. Output Predictor

Based on multi-scale feature fusion, we predict the scales and locations of objects using both score and location maps (see Figure 1(c)), which are defined as follows:

- The score map corresponds to confidence score of the object region. Similar to the confidence map in FCN [19], each pixel of the score map is a scalar between 0 to 1 representing the confidence belonging to an object region.
- The location map describes the location of object by using four distance channels $G = (l, t, r, b)$. The channels denote the distances from the current pixel i to the *left*, *top*, *right*, and *bottom* edges of the bounding box respectively. Then we can directly predict the object box by four distance channels. Specifically, for each point in the score map, four distance channels predict the distances to the above four edges of the bounding box.

2.3.1 Foreground Attention.

In general, based on both score and location maps, we can estimate the bounding boxes of the objects in the image. However, the estimated bounding boxes only rely on the global appearance of the object. That is, little local appearance of the object is taken into consideration, resulting in less discriminative foreground representation. To improve localization accuracy, we introduce the Foreground Attention (FA) module to consider both global and local appearance representation of the objects.

In practice, we use four corner maps (*top-left*, *top-right*, *bottom-left* and *bottom-right*) to denote different corner positions within the object region, as shown in Figure 3. Similar

to score map, each pixel of the corner map is also a scalar between 0 to 1 representing the confidence belonging to a corresponding position in the object region. The corner is set as 1/9 the size of the whole object. Specifically, as illustrated in Figure 1(c), we first use a threshold filter to remove the candidate bounding boxes with low confidence pixels, *i.e.*, $c_i < \mu$. c_i is the confidence value of pixel i in the predicted score map, and μ denotes the confidence threshold. Then, the Non-Maximum Suppression (NMS) operation is applied to remove redundant candidate bounding boxes and choose the top ones with higher confidence. Finally, a corner voting filter is designed to determine whether the selected bounding boxes should be retained. Specifically, we calculate the number of reliable corners $\mathbb{N}(b_k)$ in the k -th candidate bounding box b_k by

$$\mathbb{N}(b_k) = \sum_{s=1}^4 \mathbb{I}(\tau(\mathcal{C}_s) > \varepsilon), \quad (3)$$

where $\tau(\mathcal{C}_s)$ denotes the average confidence of the corner region \mathcal{C}_s . ε indicates the threshold of mean confidence $\tau(\mathcal{C}_s)$ to determine the reliable corner. $\mathbb{I}(\cdot) = 1$ if its argument is true, and 0 otherwise. We only keep the bounding box b_k if the number of reliable corners is larger than the threshold κ , *i.e.*, $\mathbb{N}(b_k) > \kappa$.

2.4. Loss function

To train the proposed network, We optimize the location map and score map, as well as both foreground and background attentions simultaneously. The overall loss function is defined as

$$\mathcal{L} = \mathcal{L}_{\text{loc}} + \lambda_{\text{sco}} \mathcal{L}_{\text{sco}} + \lambda_{\text{FA}} \mathcal{L}_{\text{FA}} + \lambda_{\text{BA}} \mathcal{L}_{\text{BA}}, \quad (4)$$

where \mathcal{L}_{loc} , \mathcal{L}_{sco} , \mathcal{L}_{FA} , and \mathcal{L}_{BA} are loss terms for the location map, score map, foreground attention, and background attention, respectively. The parameter λ_{sco} , λ_{FA} , and λ_{BA} are used to balance these terms. In the following, we explain these loss terms in detail.

2.4.1 Loss of Location Map.

To achieve scale-invariance, the IoU loss [33] is adopted to evaluate the difference between the predicted bounding box and the ground truth of bounding box. The loss of location map is defined as:

$$\mathcal{L}_{\text{loc}} = \text{IoU}(G, G^*), \quad (5)$$

where $G = (l, t, r, b)$ and $G^* = (l^*, t^*, r^*, b^*)$ are the estimated and ground-truth bounding box of the object. The function $\text{IoU}(\cdot)$ calculates the intersection-over-union (IoU) score between G and G^* .

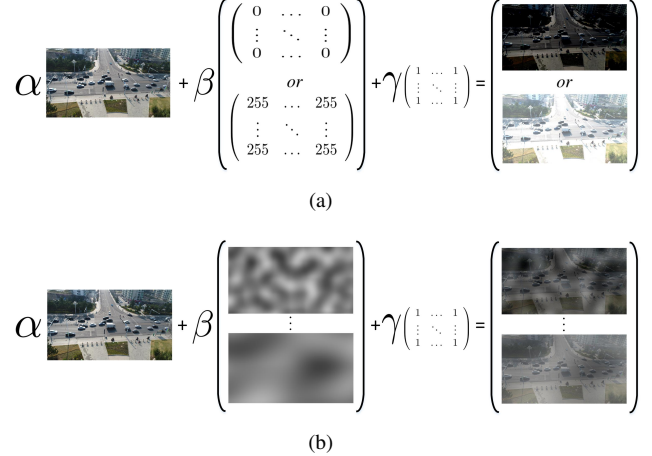


Figure 4. Illustration of data augmentation including (a) BNoise (Brightness noises to imitate sunny or night scenes) and (b) PNoise (Perlin noises to imitate cloudy and foggy scenes).

2.4.2 Loss of Score Map.

Similar to image segmentation [34], we use the Dice loss to deal with the imbalance problem of positive and negative pixels in the score map. It calculates the errors between the predicted score map and ground-truth map. The loss is calculated as

$$\mathcal{L}_{\text{sco}} = 1 - \frac{2 \cdot \sum_{i=1}^N (c_i c_i^*)}{\sum_{i=1}^N (c_i) + \sum_{i=1}^N (c_i^*)}, \quad (6)$$

where the sums run over the all N pixels of the score map. c_i^* and c_i are the confidence values of pixel i in the ground-truth and predicted maps respectively.

2.4.3 Loss of Background Attention.

Similar to classification algorithms, we use the cross-entropy loss \mathcal{L}_{BA} to guide background attention based on the binary classification, *i.e.*,

$$\mathcal{L}_{\text{BA}} = \begin{cases} -\log(p) & \text{if } y = 1, \\ -\log(1-p) & \text{otherwise,} \end{cases} \quad (7)$$

where $y \in \{\pm 1\}$ denotes the ground-truth category (*i.e.*, foreground or background), $p \in [0, 1]$ is the estimated probability for the category with label $y = 1$.

2.4.4 Loss of Foreground Attention.

Similar to the score map, to deal with the imbalance problem of positive and negative pixels in the feature maps, we use the Dice loss to guide the foreground attention for the four corner maps.

2.5. Data Augmentation for Drones

Data augmentation is important in deep network training based on limited training data. Since the data is captured from a very high altitude by the drone, it is susceptible to the influence of different illumination conditions, *e.g.*, sunny, night, cloudy and foggy. Therefore, we develop a new data augmentation strategy for drones.

As we know, sunny or night scenes correspond to the brightness of the image, therefore we synthesize these scenes via changing the whole contrast of the image (denoted as *BNoise*). On the other hand, since convincing representations of clouds and water can be created in pixel-level [21], we use Perlin noise [22] to imitate cloudy and foggy scenes (denoted as *PNoise*). Inspired by the image blending algorithm [30], the data augmentation model is defined as

$$\Phi(i) = \alpha I(i) + \beta M^*(i) + \gamma, \quad (8)$$

where $\Phi(i)$ is the transformed value of the pixel i in image. α and β denote the weight of the pixel of original image $I(p)$ and noise map $M^*(i)$ respectively. The asterisk $*$ denotes different kinds of noise maps, *i.e.*, *BNoise* $M^b(i)$ and *PNoise* $M^p(i)$. We have $\alpha = 1 - \beta$ to control the contrast of the image. The perturbation factor γ is used to revise the brightness. We set different factors α and γ for each image in the training phase.

As shown in Figure 4(a), we employ white and black maps to synthesize sunny or night images. On the other hand, we use Perlin noise [22] to generate noise maps in Figure 4(b), and then revise the brightness via disturbance factor γ to synthesize cloudy and foggy images. For each training image, we first resize it using random scale factors (x0.5, x1, x2 and x3). Then, we introduce both noise maps into the image to imitate the challenging scenes (*i.e.*, sunny, night, cloudy, and foggy). Finally, we select positive and negative images by random cropping on the blending images, and transform the selected images to 512×512 size via zooming and padding.

3. Experiment

The proposed method is implemented by Tensorflow r1.8¹. We will release the source codes of our method upon the acceptance of the paper. We evaluate our method on two drone based datasets: UAVDT [4] and CARPK [10]. We also evaluate our method on the PUCPR+ dataset [10] because the dataset is collected from the 10th floor of a building and similar to drone view images to a certain degree. In this section, we first describe implementation details. Then, we compare our GANet with the state-of-the-art methods, *i.e.*, Faster R-CNN [25], RON [12], SSD [18], R-FCN [2], CADNet [5], One-Look Regression [20], IEP [29],

¹<https://www.tensorflow.org/>

YOLO9000 [23], LPN [10], RetinaNet [17], YOLOv3 [24], IoUNet [8], and SA+CF+CRT [15]. More visual examples are shown in Figure 5. In addition, the ablation study is carried out to evaluate the effectiveness of each component in our network.

3.0.1 Implementation Details

Due to the shortage of computational resources, we train GANet using the VGG-16 and ResNet-50 backbone with the input size 512×512 . All the experiments are carried out on the machine with NVIDIA Titan Xp GPU and Intel(R) Xeon(R) E5-1603v4@2.80GHz CPU. For fair evaluation, we generate the same top 200 detection bounding boxes for the UAVDT and CARPK datasets and 400 detection bounding boxes for the PUCPR+ dataset based on the detection confidence. Note that the detection confidence is calculated by summarizing the value of each pixel in the score map. To output the count of objects in each image, we calculate the number of detection with the detection confidence larger than 0.5. We fine-tune the resulting model using the Adam Optimizer. An exponential decay learning rate is used in the training phrase, *i.e.*, its initial value is 0.0001 and decays every 10,000 iterations with the decay rate 0.94. The batch size is set as 10. In the loss function (4), we set the balancing factors as $\lambda_{sco} = 0.01$, $\lambda_{FA} = 0.0025$, $\lambda_{BA} = 0.001$ empirically. In the FA module, the confidence threshold μ is set as 0.8, and the threshold ε in (3) is set as 0.3 empirically. The Non-Maximum Suppression (NMS) operation is conducted with a threshold 0.2. In the data argumentation model (8), we set the balancing weights as $\alpha = \{0.1, 0.3, 0.5, 0.7, 0.8, 0.9, 1.0\}$ and $\gamma = [-20, 20]$.

3.0.2 Metrics.

To evaluate detection algorithms on the UAVDT dataset [4], we compute the Average Precision (AP@0.7) score based on [7, 6]. That is, the hit/miss threshold of the overlap between detection and ground-truth bounding boxes is set to 0.7. In terms of CARPK [10] and PUCPR+ [10], we report the detection score under two hit/miss thresholds, *i.e.*, AP@0.5 and AP@0.7. To evaluate the counting results, similar to [10], we use two object counting metrics including Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE).

3.1. Quantitative Evaluation

3.1.1 Evaluation on UAVDT.

The UAVDT dataset [4] consists of 100 video sequences with approximate 80,000 frames, which are collected from various scenes. Moreover, the objects are annotated by bounding boxes as well as several attributes (*e.g.*, weather condition, flying altitude, and camera view). Note that we only use the

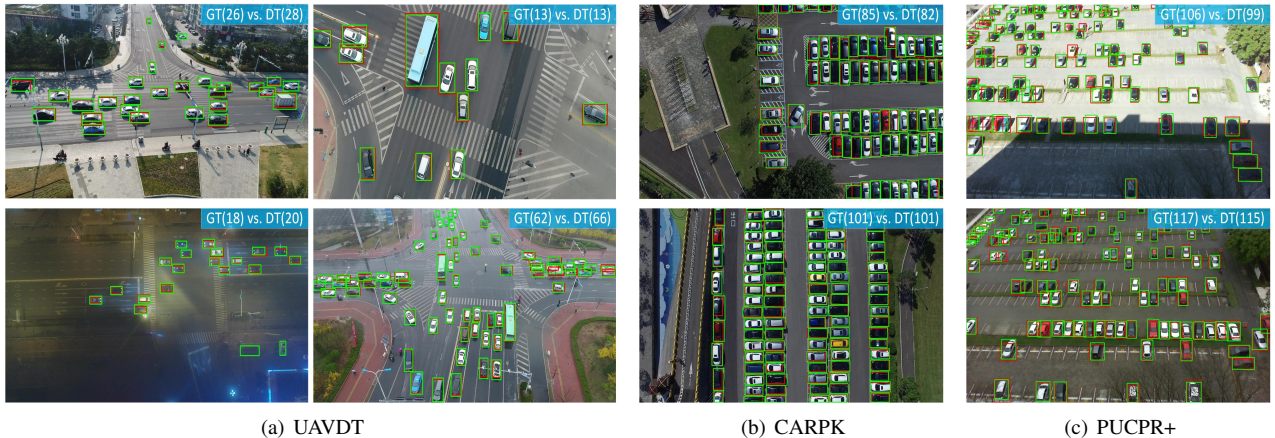


Figure 5. Visual examples of GANet with VGG16 backbone. The ground-truth and predicted detection bounding boxes are highlighted in red and green rectangles, respectively. The blue mask in the top-right corner indicates the comparison between the ground-truth (GT) and estimated detection (DT) counts.

Table 1. Comparison on the UAVDT dataset.

Method	Backbone	MAE↓	RMSE↓	AP@0.7[%]↑
YOLO9000	DarkNet-19	12.59	16.73	7.6
YOLOv3	DarkNet-53	11.58	21.50	20.3
RON	VGG-16	-	-	21.6
Faster R-CNN	VGG-16	-	-	22.3
SSD	VGG-16	-	-	33.6
CADNet	VGG-16	-	-	43.6
Ours	VGG-16	5.10	8.10	46.8
SA+CF+CRT	ResNet-101	7.67	10.95	27.8
R-FCN	ResNet-50	-	-	34.4
Ours	ResNet-50	5.09	8.16	47.2

Table 2. Comparison on the CARPK dataset.

Method	MAE↓	RMSE↓	AP@0.5[%]↑	AP@0.7[%]↑
One-Look Regression	59.46	66.84	-	-
IEP	51.83	-	-	-
Faster R-CNN	47.45	57.39	-	-
YOLO9000	38.59	43.18	20.9	3.7
SSD	37.33	42.32	68.7	25.9
LPN	23.80	36.79	-	-
RetinaNet	16.62	22.30	-	-
YOLOv3	7.92	11.08	85.3	47.0
IoUNet	6.77	8.52	-	-
SA+CF+CRT	5.42	7.38	89.8	61.4
Ours (VGG-16)	4.80	6.94	90.2	73.6
Ours (ResNet-50)	4.61	6.55	90.1	74.9

subset of UAVDT dataset for object detection in our experiment. As presented in Table 1, we can conclude that our GANet performs the best among all the compared detection methods in terms of both the VGG-16 and ResNet-50 backbones. Specifically, GANet surpasses YOLO9000, YOLOv3,

Table 3. Comparison on the PUCPR+ dataset.

Method	MAE↓	RMSE↓	AP@0.5[%]↑	AP@0.7[%]↑
SSD	119.24	132.22	32.6	7.1
Faster R-CNN	111.40	149.35	-	-
YOLO9000	97.96	133.25	12.3	4.5
RetinaNet	24.58	33.12	-	-
LPN	23.80	36.79	-	-
One-Look Regression	21.88	36.73	-	-
IEP	15.17	-	-	-
IoUNet	7.16	12.00	-	-
YOLOv3	5.24	7.14	95.0	45.4
SA+CF+CRT	3.92	5.06	92.9	55.4
Ours (VGG-16)	3.68	5.47	91.3	67.0
Ours (ResNet-50)	3.28	4.96	91.4	65.5

RON, Faster R-CNN, SSD, CADNet, SA+CF+CRT and R-FCN by 39.2%, 26.3%, 25.2%, 24.5%, 13.2%, 3.2%, 19.4% and 12.8% AP scores, respectively. Moreover, our method achieves better counting accuracy than SA+CF+CRT with the more complex ResNet-101 backbone, *i.e.*, 5.09 MAE score and 8.16 RMSE score. It demonstrates that the effectiveness of our method in object detection in drone based scenes.

3.1.2 Evaluation on CARPK.

The CARPK dataset [10] provides the largest-scale drone view parking lot dataset in unconstrained scenes, which is collected in various scenes for 4 different parking lots. It contains approximately 90,000 cars in total with the view of drone. We compare our method with state-of-the-art algorithms in Table 2. The results show that our approach achieves the best MAE, RMSE and AP scores. It is worth mentioning that we obtain much better AP@0.7 score (*i.e.*, 74.9 vs. 61.4). This is attributed to the proposed attention

Table 4. Comparison of variants of GANet on the UAVDT dataset.

Method	AP	AP _{day}	AP _{night}	AP _{fog}	AP _{low}	AP _{med}	AP _{high}	AP _{front}	AP _{side}	AP _{bird}
GANet	0.3908	0.4779	0.5513	0.1509	0.5505	0.4616	0.1227	0.4478	0.5111	0.1981
GANet+BPNoise	0.4181	0.4940	0.5581	0.2027	0.5565	0.4867	0.1665	0.4618	0.5219	0.2533
GANet+FA	0.4207	0.5006	0.5878	0.1890	0.5935	0.4834	0.1431	0.4595	0.5462	0.2439
GANet+BA	0.4353	0.5041	0.5743	0.2401	0.5908	0.4812	0.1996	0.4655	0.5451	0.2947
GANet+BPNoise+FA	0.4411	0.5272	0.5819	0.2139	0.5900	0.5146	0.1751	0.4805	0.5618	0.2715
GANet+BPNoise+BA	0.4576	0.5049	0.5779	0.3068	0.5815	0.4923	0.2695	0.4719	0.5309	0.3640
GANet+BPNoise+FA+BA	0.4679	0.5240	0.5841	0.3084	0.5820	0.5206	0.2624	0.4852	0.5435	0.3603

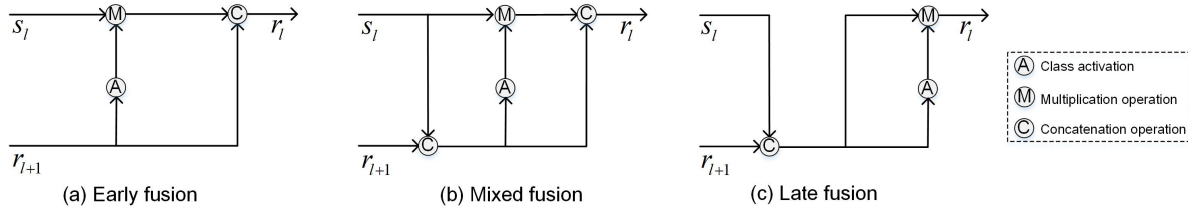
Figure 6. Different fusion strategies of multi-scale feature maps. s_l denotes the low-level features with rich texture details, r_{l+1} and r_l denote the high-level features with strong semantic information.

Table 5. Influence of data augmentation.

Method	AP	AP _{day}	AP _{night}	AP _{fog}
GANet	0.3908	0.4779	0.5513	0.1509
GANet+BNoise	0.4034	0.4928	0.5686	0.1579
GANet+PNoise	0.4063	0.4798	0.5263	0.2118
GANet+BPNoise	0.4181	0.4940	0.5581	0.2027

Table 6. Influence of background attention.

Method	AP	AP _{front}	AP _{side}	AP _{bird}
GANet+BPNoise	0.4181	0.4618	0.5219	0.2533
GANet+BPNoise+LF	0.4457	0.4667	0.5301	0.3294
GANet+BPNoise+MF	0.4530	0.4699	0.5338	0.3495
GANet+BPNoise+EF	0.4576	0.4719	0.5309	0.3640
GANet+BPNoise+FPN	0.3985	0.4378	0.4943	0.2480
GANet+BPNoise+GC	0.4343	0.4681	0.5374	0.2919
GANet+BPNoise+SE	0.4442	0.4723	0.5347	0.3142
GANet+BPNoise+BA	0.4576	0.4719	0.5309	0.3640

modules to locate the objects more accurately.

3.1.3 Evaluation on PUCPR+.

The PUCPR+ dataset [10] is the subset of PKLot [3], which is annotated with nearly 17,000 cars in total. It shares the similar high altitude attribute to drone based scenes, but the camera sensors are fixed and set in the same place. As presented in Table 3, our method performs the best in terms of MAE and RMSE scores. YOLOv3 [24] achieves the best AP score at 0.5 hit/miss threshold, but inferior AP@0.7 score than that of our method. We speculate that YOLOv3 lack of global appearance representation of objects to achieve accurate localization.

3.2. Ablation Study

We perform analyses on the effect of the important modules in our method on the detection performance. Specifically, we study the influence of data augmentation, semantic discriminative attention, and corner attention. We select the UAVDT dataset [4] to conduct the experiment because it provides various attributes in terms of altitude, illumination

and camera-view for comprehensive evaluation.

3.2.1 Effectiveness of Data Augmentation.

As discussed above, the data augmentation strategy is used to increase the difficult samples affected by various illumination attributes in the UAVDT dataset [4] such as *daylight*, *night* and *fog*. We compare different variants of GANet with different data augmentation, denoted as GANet+BNoise, GANet+PNoise and GANet+BPNoise. Notably, BNoise denotes the brightness noise, PNoise denotes the Perlin noise, and BPNoise denotes both. As shown in Table 5, the performance of GANet+BNoise is slightly higher than that of GANet. GANet+PNoise achieves much better AP score in terms of foggy scenes compared to GANet (0.2118 vs. 0.1509), which demonstrates the effectiveness of the introduced Perlin noise. If we perform the full data augmentation strategy in our training samples, the overall performance will increase by 2%.

Table 7. Influence of foreground attention.

Method	κ	AP	AP _{low}	AP _{med}	AP _{high}
GANet+BPNoise	-	0.4181	0.5565	0.4867	0.1656
	0	0.4271	0.5729	0.4978	0.1698
	1	0.4411	0.5900	0.5146	0.1751
GANet+BPNoise+FA	2	0.4391	0.5869	0.5130	0.1736
	3	0.4372	0.5817	0.5128	0.1718
	4	0.4347	0.5764	0.5116	0.1699

3.2.2 Effectiveness of Background Attention.

Different from the previous unsupervised attention modules, our Background Attention (BA) is guided based on discrimination between the background and objects. Firstly, we study different fusion strategies of the proposed BA in Figure 6, *i.e.*, early fusion (EF), mixed fusion (MF) and late fusion (LF). The results presented in Table 6 show the early fusion strategy (*i.e.*, GANet+BPNoise+EF) achieves the best performance. Secondly, we also compare BA with several previous channel-wise attention modules including SE block [11] and GC block [1]. For a fair comparison, we use the same early fusion strategy in Figure 6(a). Compared to the baseline FPN fusion strategy using lateral connection [16], all the attention modules can improve the performance by learning the weights of different channels of feature maps. However, our BA module can learn additional discriminative information of background, resulting in the best AP score in the drone based scenes under different camera views (*i.e.*, *front-view*, *side-view* and *bird-view*).

3.2.3 Effectiveness of Foreground Attention.

We enumerate the threshold for Foreground Attention (FA) κ in (3), *i.e.*, $\kappa = \{0, 1, 2, 3, 4\}$, to study its influence on the accuracy. As shown in Table 7, we can conclude that GANet with the FA module achieves the best AP score 0.4411 when the threshold $\kappa = 1$. If we remove FA, the detection performance will decrease to 0.4181. It shows the effectiveness of the FA module.

3.2.4 Variants of GANet.

In Table 4, we compare various variants of GANet that combine several components in the network. Using data argumentation strategy can improve the performance considerably in all the attributes. Either BA or FA can improve the performance by 3 ~ 4%. Moreover, the proposed method using both attentions and data argumentation strategy can boost the performance by approximate 8% improvement in AP score compared to the baseline GANet method.

4. Conclusion

In the paper, we propose a novel guided attention network to deal with object detection and counting in drone based scenes. Specifically, we introduce both background and foreground attention modules to not only learn background discriminative representation but also consider local appearance of the object, resulting in better accuracy. The experiments on three challenging datasets demonstrate the effectiveness of our method. We plan to expand our method to multi-class object detection and counting for future work.

References

- [1] Y. Cao, J. Xu, S. Lin, F. Wei, and H. Hu. Gcnet: Non-local networks meet squeeze-excitation networks and beyond. *CoRR*, abs/1904.11492, 2019.
- [2] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: object detection via region-based fully convolutional networks. In *NeurIPS*, pages 379–387, 2016.
- [3] P. R. L. de Almeida, L. S. Oliveira, A. S. B. Jr., E. J. da Silva Junior, and A. L. Koerich. Pklot - A robust dataset for parking lot classification. *Expert Syst. Appl.*, 42(11):4937–4949, 2015.
- [4] D. Du, Y. Qi, H. Yu, Y. Yang, K. Duan, G. Li, W. Zhang, Q. Huang, and Q. Tian. The unmanned aerial vehicle benchmark: Object detection and tracking. In *ECCV*, pages 375–391, 2018.
- [5] K. Duan, D. Du, H. Qi, and Q. Huang. Detecting small objects using a channel-aware deconvolutional network. *TCSVT*, 2019.
- [6] M. Everingham, S. M. A. Eslami, L. J. V. Gool, C. K. I. Williams, J. M. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *IJCV*, 111(1):98–136, 2015.
- [7] A. Geiger, P. Lenz, and R. Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. In *CVPR*, pages 3354–3361, 2012.
- [8] E. Goldman, R. Herzig, A. Eisenschlat, O. Ratzon, I. Levi, J. Goldberger, and T. Hassner. Precise detection in densely packed scenes. *CoRR*, abs/1904.00853, 2019.
- [9] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.
- [10] M. Hsieh, Y. Lin, and W. H. Hsu. Drone-based object counting by spatially regularized regional proposal network. In *ICCV*, pages 4165–4173, 2017.
- [11] J. Hu, L. Shen, and G. Sun. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141, 2018.
- [12] T. Kong, F. Sun, A. Yao, H. Liu, M. Lu, and Y. Chen. RON: reverse connection with objectness prior networks for object detection. In *CVPR*, pages 5244–5252, 2017.
- [13] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012.
- [14] H. Law and J. Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, pages 765–781, 2018.

- [15] W. Li, H. Li, Q. Wu, X. Chen, and K. N. Ngan. Simultaneously detecting and counting dense vehicles from drone images. *TIE*, 66(12):9651–9662, 2019.
- [16] T. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie. Feature pyramid networks for object detection. In *CVPR*, pages 936–944, 2017.
- [17] T. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, pages 2999–3007, 2017.
- [18] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. E. Reed, C. Fu, and A. C. Berg. SSD: single shot multibox detector. In *ECCV*, 2016.
- [19] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, pages 3431–3440, 2015.
- [20] T. N. Mundhenk, G. Konjevod, W. A. Sakla, and K. Boakye. A large contextual dataset for classification, detection and counting of cars with deep learning. In *ECCV*, pages 785–800, 2016.
- [21] K. Perlin. An image synthesizer. In *SIGGRAPH*, pages 287–296, 1985.
- [22] K. Perlin. Improving noise. *TOG*, 21(3):681–682, 2002.
- [23] J. Redmon and A. Farhadi. YOLO9000: better, faster, stronger. In *CVPR*, pages 6517–6525, 2017.
- [24] J. Redmon and A. Farhadi. Yolov3: An incremental improvement. *CoRR*, abs/1804.02767, 2018.
- [25] S. Ren, K. He, R. B. Girshick, and J. Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In *NeurIPS*, pages 91–99, 2015.
- [26] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241, 2015.
- [27] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. *CoRR*, abs/1612.06851, 2016.
- [28] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [29] T. Stahl, S. L. Pinteá, and J. C. van Gemert. Divide and count: Generic object counting by image divisions. *TIP*, 28(2):1035–1044, 2019.
- [30] R. Szeliski. *Computer vision: algorithms and applications*. Springer Science & Business Media, 2010.
- [31] X. Wang, R. B. Girshick, A. Gupta, and K. He. Non-local neural networks. In *CVPR*, pages 7794–7803, 2018.
- [32] Z. Yang, S. Liu, H. Hu, L. Wang, and S. Lin. Repoints: Point set representation for object detection. *CoRR*, abs/1904.11490, 2019.
- [33] J. Yu, Y. Jiang, Z. Wang, Z. Cao, and T. S. Huang. Unitbox: An advanced object detection network. In *ACM MM*, 2016.
- [34] J. Zhang, X. Shen, T. Zhuo, and H. Zhou. Brain tumor segmentation based on refined fully convolutional neural networks with a hierarchical dice loss. *CoRR*, abs/1712.09093, 2017.
- [35] T. Zhou, P. Krähenbühl, M. Aubry, Q. Huang, and A. A. Efros. Learning dense correspondence via 3d-guided cycle consistency. In *CVPR*, pages 117–126, 2016.
- [36] X. Zhou, D. Wang, and P. Krähenbühl. Objects as points. *CoRR*, abs/1904.07850, 2019.