# FOCUS-AND-DETECT: A SMALL OBJECT DETECTION FRAMEWORK FOR AERIAL IMAGES *

**Onur Can Koyun**
İstanbul Technical University
İstanbul, Turkey
okoyun@itu.edu.tr

**Reyhan Kevser Keser**
Istanbul Technical University
İstanbul, Turkey
keserr@itu.edu.tr

**İbrahim Batuhan Akkaya**
Aselsan
Ankara, Turkey
ibakkaya@aselsan.com.tr

**Behçet Uğur Töreyin**
Istanbul Technical University
İstanbul, Turkey
toreyin@itu.edu.tr

## ABSTRACT

Despite recent advances, object detection in aerial images is still a challenging task. Specific problems in aerial images makes the detection problem harder, such as small objects, densely packed objects, objects in different sizes and with different orientations. To address small object detection problem, we propose a two-stage object detection framework called "Focus-and-Detect". The first stage which consists of an object detector network supervised by a Gaussian Mixture Model, generates clusters of objects constituting the focused regions. The second stage, which is also an object detector network, predicts objects within the focal regions. Incomplete Box Suppression (IBS) method is also proposed to overcome the truncation effect of region search approach. Results indicate that the proposed two-stage framework achieves an AP score of 42.06 on VisDrone validation dataset, surpassing all other state-of-the-art small object detection methods reported in the literature, to the best of authors' knowledge.

*Keywords* Object detection · Small object detection · Region search · Aerial images

## 1 Introduction

Object detection is a computer vision task which consists of two sub-tasks, namely, object localization and classification. It is one of the fundamental problems, since many other tasks rely on it, such as image captioning, object tracking, instance segmentation and scene understanding [1]. Thus, it has been studied for a long time. With the progress of deep learning based methods, handcrafted feature based methods, such as HOG [2] and SIFT [3], have become obsolete. SIFT and HOG features are low-level features which cannot be utilized as hierarchical layer-wise representations while the deep models are able to represent the data as hierarchical combination of abstract representations. Nevertheless, recent methods are getting more complex day by day thanks to the development on hardware capabilities. In [4], deep learning based methods are defined as a combination of various components. In general, detection networks consist of backbone, neck and head. In this context, backbone model is the network that extracts features for the detection task, head is the actual detection model that predicts both bounding boxes and classes, neck is placed between backbone and head networks and fuses feature maps from different stages of backbone model. There are different approaches for detection heads, such as one-stage detection and two-stage detection models. One-stage detection models do not include a region proposal layer [5] in the head model and run detection directly over a dense sampling of locations. On
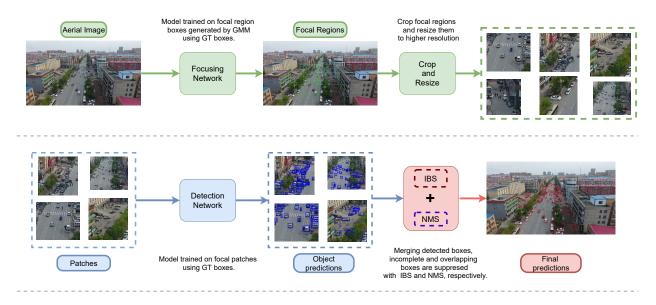
Figure 1: Focus&Detect (F&D) framework. F&D consists of two components : (1) Focus network; (2) Detection network. While focus network detects the possible object containing regions (Focal Regions), detection network detects objects in these focal regions. Final predictions are generated by merging the predictions of focal regions. NMS and IBS methods are applied to eliminate overlapping and truncated boxes. Both detectors are trained with supervision. Focus network utilizes cluster coordinates generated by a Gaussian Mixture Model as supervision signal. On the other hand, detection network utilizes object ground truth bounding boxes in each respective focal region.

the other hand two-stage models utilize region proposal network to extract object regions which are used for bounding box regression and classification.

Aerial object detection, which can be categorized as a case of the general small object detection problem, is an emerging field with recent advances. Although, it has a wide range of applications, such as surveillance, precision agriculture, military monitoring and urban management [6, 7], it is one of the most challenging computer vision tasks. Earlier, several studies proposed adapting methods established for natural images to aerial images [8, 9]. However, various difficulties arose due to such an approach [10]. First, in aerial images, orientation and aspect ratios may substantially differ from natural images. Second, scale variation is much severe in aerial images for both intra-class and inter-class samples [11]. As an example, [12] reports the statistics for 'car' class in the MS COCO and the VisDrone [13] datasets. It is indicated that in the VisDrone dataset, the variance of the sizes of 'car' objects is almost five times larger than that of the MS COCO dataset. Third, objects in aerial images are small and densely placed. For example, up to 902 objects may exist in a single image in VisDrone Detection dataset [14]. Moreover, class imbalance problem exists in aerial images [14], where it makes the small object detection problem even harder for classes with small number of samples. Hence, dedicated approaches addressing the aforementioned problems are required for the small object detection task.

Region search is a powerful method for small object detection, which aims to find and focus on regions that potentially include objects [15, 16]. Since aerial images consist of dense and small objects, we focus on region search for aerial object detection problem, in this paper. For this purpose, we propose a framework consisting of two stages, namely the focus and the detection stages. In the first stage, regions to be focused are determined by a detector which is supervised by a Gaussian Mixture Model. The second stage, fed by these regions which are mainly clusters of objects, predicts objects within these regions. While merging the predictions on these regions, NMS and the proposed IBS methods are utilized to eliminate overlapping and truncated bounding boxes.

Our contributions can be listed as follows:

- We propose a framework, namely, 'Focus&Detect' for small object detection in aerial images which is based on region search

- We propose a method to generate object clusters using Gaussian mixture model, where the generated clusters are scale normalized.

- We also propose the 'Incomplete Box Suppression' (IBS) approach to suppress incomplete boxes caused by overlapping focal regions.

- Our proposed method achieves 42.06 AP score on the VisDrone validation set and 54.16 AP@70 score on UAVDT test set. To the best of our knowledge, our method outperforms state-of-the-art methods for small object detection, reported in the literature on VisDrone dataset.

## 2 Related Work

In this section we briefly review related work in the directions of object detection, small object detection and object detection in aerial images.

### 2.1 Object Detection

Recent methods in object detection literature typically employ powerful models for backbone, such as ResNet [17], Hourglass [18] and ResNeXt [19]. Feature pyramid network [20] based architectures are the main choice for neck model. There are multi-stage head models in the literature such as Faster-RCNN [5], Mask-RCNN [21], Cascade-RCNN [22]. Faster R-CNN generates proposals by the region proposal network (RPN). Mask R-CNN extends Faster R-CNN to perform detection and segmentation tasks simultaneously. On the other hand, YOLOv3 [23], SSD [24], GFL [25] and RetinaNet [26] are examples of single-stage detectors. Single-stage detectors omit the proposal stage and make detection on the dense sample of locations.

Recent two-stage methods include RPN [5] based modules to generate region proposals. Anchor is an alternative concept in object detection literature standing for pre-defined bounding boxes that are to be matched with ground-truth bounding boxes of objects. Some studies propose anchor-free approaches to avoid the computational cost of using anchors [27, 28]. A key component for deep learning based object detection methods is the loss function. It mainly consists of two terms corresponding to the regression and the classification losses [5]. Regression and classification branches are the final modules of an object detector, which predict the localization and classification of objects, respectively. Finally, Non-Maximum Suppression (NMS) and its variations [29, 30] are instrumental in the workflow of an object detector, since object detectors typically generate lots of redundant predictions and NMS is used for reducing redundancy. With the improvements on these components, great progress has been made on generic object detection [31, 32], whereas small object detection still needs improvements to obtain satisfying detection performances [4]. Adapting an object detector which performs adequately well for medium or large objects, cannot yield sufficient performances on small objects, whose areas are less than $32 \times 32$ pixels as stated in [33]. For instance, it is indicated that recently proposed DETR's [34] detection performance is not at the desired level for small objects, whereas DETR performs better than Faster R-CNN [5] on MS COCO [33] dataset.

### 2.2 Small Object Detection

Small object detection task is an important computer vision problem with applications in various fields, such as autonomous driving, UAV-based imaging, and surveillance. Even though, it is a crucial tool for different computer vision tasks in numerous fields, performances of the current methods are not at the desired level. Still, most of the object detection methods struggle with small objects due to the issues, such as the inadequacy of information raised from the area covered by small objects on the image, high possibility of location for small objects and being adapted for medium and large objects [1].

To solve the inadequacy of information problem for small object detection, [35] uses super-resolution techniques to improve the performance of Faster R-CNN, whereas [36] utilizes a super-resolution GAN on remote sensing images.

Increasing the resolution of input yield better performance for small objects. Thus, some trivial methods are also proposed such as using an image pyramid as the input to improve the performance of detecting small faces [37]. However these methods are not scalable efficiently.

In [38], a two-stream network is proposed which utilizes multi-scale representation as well as the attention mechanism. Another study [39] uses contextual information besides the multi-scale representation obtained from SSD model. Pan et. al.[40] proposed a multi-scale feature fusing scheme to improve small object detection on SSD model. In another study [41], SSD model is revised with feature fusion and dilated convolutions increasing the detection performance.

On the other hand, some methods are focused on improving the region-proposal stage using various techniques such as advancing anchors [42] and increasing samples of small objects [43]. In [44], authors proposed a hybrid model which uses both region-proposal and dense detection heads to increase performance.

## 2.3 Object Detection in Aerial Images

Aerial images constitute one of the most difficult cases for object detection as they mostly comprise small objects, large difference between number of samples of different classes and the high scale variance on both of inter-class and intra-class. To alleviate these difficulties, numerous methods are proposed previously. For example, an adaptive augmentation method is proposed for the class imbalance problem in [45], which is called AdaResampling. In [46], a hard chip mining method is proposed as data augmentation on aerial images. Moreover, [11] proposes an improvement on obtaining multi-scale features in order to reduce the effect of scale variance for object detection.

Since aerial images mostly consist of small and dense objects, some methods focus on improving region search [46, 47, 48, 49, 16, 50, 15, 51]. For example, [48] proposes tiling based method to detect pedestrians and vehicles in aerial images in real-time. In [16], difficult cluster regions are determined using mean shift algorithm to feed the object detector. [50] proposes three augmentation methods for cropping based approach which are mosaic augmentation, adaptive cropping and mask resampling. In [12], an adaptive image cropping method based on FPN [20] is proposed to solve scale challenges in aerial images.[47] constructs density maps to determine regions to be cropped. Then an object detector is fed by these crops as well as the whole image. [15] utilizes clustering to obtain image crops. Before feeding object detector with these crops, this method re-scales them by determining the proper scales for the objects, to avoid the degradation on performance. Contrary to [47, 15], our method utilizes predicted regions only, and does not employ detection on the whole image. Furthermore, [47] uses density maps to generate regions of interest, which does not directly normalize the object scales. [15] utilizes a sub-network to predict scales of detected clusters which means additional computation. On the other hand, Gaussian mixture model provides scale normalization across predicted regions without additional computation as resizing predicted regions to a fixed size, yields a shift in the mean of each mixture component and resulting normalization of bounding boxes.

Different from the previous studies, we propose to use Gaussian Mixture Model (GMM) for region search. Moreover, we propose Incomplete Box Suppression (IBS) in order to suppress incomplete boxes within overlapping regions that are generated by the first detector under supervision of GMM. Figure 2 demonstrates the contribution of proposed IBS method.
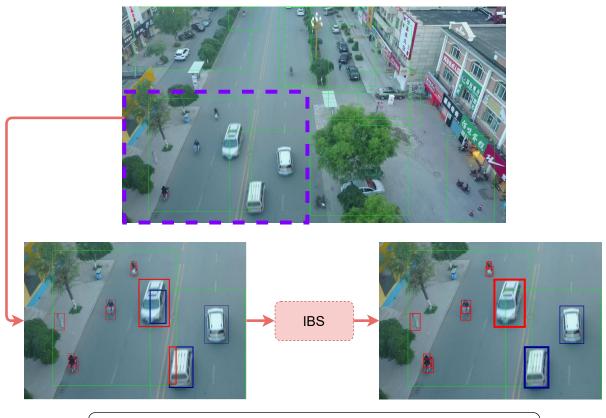
# 3 Focus-and-Detect

## 3.1 Overview

In general, object detection performance on aerial images are hindered by the small objects, changes in the perspective of objects, occlusion and truncation. Using high resolution images as input is one of the simplest solutions to the small object detection problem. Unfortunately, high-resolution images impose an unaffordable computational cost to deep neural networks. Using a focusing mechanism and increasing the resolution of the focal regions have the advantages of this simple method, but at lower computational cost. As shown in Figure 1, detection on aerial images consist of two stages: Focus network which detects focal regions consisting of cluster of objects, detection network which detects objects in focal regions. Post processing methods are applied after merging the predictions. Specifically, we proposed the Incomplete Box Suppression (IBS) mechanism to suppress incomplete boxes from overlapping focal regions. We also use standard non-max suppression (NMS) to suppress overlapping boxes after merging the predicted boxes.

## 3.2 Focus Stage

Focus stage consists of an object detection network, trained to detect focal regions. Focal regions are generated using a Gaussian Mixture Model via ground truth bounding boxes. Generalized focal loss (GFL) [25] is selected as base detection method. Backbone of the model is ResNet-50 network with deformable convolutional layers [52].

Second part of model, namely, Feature Pyramid Network (FPN) aims to exploit and refine the feature maps obtained from different stages of ResNet-50, and the last part is detection head of the model which predicts bounding boxes of focal regions. The deformable convolution layer is used in the last three-stage of the backbone.

The traditional convolutional network has limited performance on geometric transformation due to the restricted form of convolutional layers and pooling layers [53]. The traditional network architecture is not able to transfer well on the focal region detection task. Transferability of focal region features are inferior to transferability of traditional object features. In order to improve the transferability of the learned features, the deformable convolutional layers [52] has been utilized within ResNet-50, since deformable convolution can change the reception field dynamically. The proposed change leads to better representation of focal regions.

Overlapped Focal Regions yield incomplete boxes, IBS suppresses the incomplete boxes in overlapped regions.

Figure 2: Focus-and-Detect predicts focal regions comprising object clusters, and detects objects in these regions. Merging predictions of all focal regions constitute the final step. Overlapping focal regions might yield incomplete bounding box predictions that fail to fully cover the whole areas of objects. The blue box contained within the red one and the red box within the blue one in the lower left subfigure are examples of such incomplete boxes. Incomplete boxes cause wrong predictions for object class affiliations. To overcome this issue, we proposed a method called Incomplete Box Suppression (IBS), where each box prediction in a focal region is able to suppress predictions of other focal regions, yielding only the 'complete' predictions.

The performance of the overall framework mostly depends on the focus stage. Ideally, predicted focal regions must include all the object bounding boxes without any truncation. However, there might be overlapped regions and truncated objects in these regions. These issues are resolved by employing the IBS method as a post-processing stage, and presented in detail below.

### 3.2.1 Generating Ground-Truth Boxes of Focal Regions Using Gaussian Mixture Model

In object localization problem, areas of objects in the same class can be modeled with a Gaussian distribution, as object sizes do not vary much. This assumption is true for object detection datasets such as MS COCO [33] or PASCAL VOC [54]. However, in aerial image datasets such as VisDrone [14], object areas deviate from one to another depending on the angle and the altitude of the camera. Instead of a single Gaussian model, a Gaussian mixture model is a better choice whereas, contrary to a single Gaussian model, a mixture model consist of Gaussians with smaller deviations when object locations are used as input to the mixture model.

In this context, focal regions can be defined as clusters of objects which are obtained with a Gaussian mixture model that takes the location information of ground-truth (GT) boxes as input. The location information consists of a vector of bounding box distances to grid of evenly sampled points in the image as seen in Figure 3. This method yields better results compared to directly using coordinates of boxes.
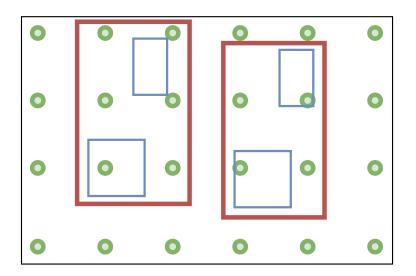
5

Figure 3: Distance vectors are defined as vectors emanating from the predefined grid of points and ending at the center of the ground truth bounding box coordinates. Distance vectors are then input to GMMs for clustering bounding boxes at each image. Blue boxes represent GT boxes of objects, red boxes represent generated focal regions.

Number of focal regions is selected depending on the number of GT boxes for the respective image. Number of focal regions ($N_f$) can be written as:

$$N_f = log_2(N_{gt}) + 2 \qquad (1)$$

where $N_{gt}$ is the number of GT boxes. Let $\vec{x}$ be a $1 \times M$ sized distance vector of the $i^{th}$ GT box in an image, and $X$ be the $N_f \times M$ sized array of feature vectors. Gaussian Mixture Model can be defined as:

$$p_{\vec{x}} = \sum_{j=1}^{N_f} \phi_j \mathcal{N}(\vec{x}|\vec{\mu_j}, \Sigma_j),$$

$$\sum_{j=1}^{N_f} \phi_j = 1 \qquad (2)$$

where $\mu_j$ and $\sigma_j$ are the mean and variance of the $j^{th}$ cluster.

Expectation maximization algorithm is used to fit the model. Once the EM algorithm has run to completion, the fitted model can be used to perform clustering on GT bounding boxes. Given the model's parameters, probability that a GT bounding box belongs to a cluster is calculated as:

$$p_{C_i|\vec{x}} = \frac{\phi_i \mathcal{N}(\vec{x}|\vec{\mu_i}, \Sigma_i)}{\sum_{j=1}^{N_f} \phi_j \mathcal{N}(\vec{x}|\vec{\mu_j}, \Sigma_j)} \qquad (3)$$

After calculation of clusters, focal regions are selected as the minimum sized box that includes all bounding boxes with a 20 pixel gap on each side in respective cluster. Because of the gap, there might be truncated objects in focal regions. Generated focal regions are used as ground truth bounding boxes for focus stage as seen in Figure 4.

### 3.3 Detection Stage

After obtaining the focal regions, a dedicated detector is utilized to perform object detection on these regions. Obtained regions are resized to a higher resolution. This approach improves performance of small object detection.

In this stage, Generalized Focal Loss (GFL) is adopted as the base detector. Backbone of the model is selected as ResNeXt-101 network with deformable convolutional layers. On the neck, Feature pyramid network (FPN) is used to improve detection performance by using features from different stages, and the last part is the detection head of the model which predicts bounding boxes of objects. The deformable convolution layer is used in the last three stages.

Figure 4: Focal region ground truth examples. For each image, we fit a Gaussian Mixture Model. From generated clusters, target bounding boxes of focal regions are obtained.

The deformable convolution yields better results than traditional convolutional layers at detecting small objects, as it is able to dynamically change its receptive field and improves the detection performance which is hindered by geometric transformations.

On detection stage, focal regions that are obtained with GMM are cropped and resized to gather a new set of data. GT bounding boxes are obtained and refined to focal region crops. Truncated GT boxes are included if at least $30\%$ of the box resides within the cropped region.

### 3.4  Post Processing

To obtain final predictions of object bounding boxes, predictions from detection stage must be merged as model outputs predictions of focal regions. The post processing steps that are applied to improve performance consist of Incomplete Box Suppression (IBS) and Non-Max Suppression (NMS).

#### 3.4.1  Incomplete Box Suppression

Models that are utilizing region search have certain problems. For instance, merging detections of target regions might be difficult, as there might be overlapped regions and truncated objects. This problem yields multiple bounding box predictions on same object. Because of the truncation, predicted bounding boxes are not fully overlapped. Thus, non-max suppression is not able to suppress these kind of false predictions. However, these predictions decrease the AP score.

In general, non-max suppression is used to eliminate highly overlapping boxes. It works well enough for traditional object detection problem. However, in most of the region search approaches, there is a final step which is merging the predictions of target regions. This creates a new problem. Overlapped regions and truncated objects in these regions which lowers the overall performance as detector might predict a full version of bounding box, and a truncated version of the bounding box for the same object as shown in Figure 5. Generally, intersection over union of these bounding boxes are small. Thus, they are able to escape from NMS. Truncated objects are also a problem on their own. False class predictions are common for truncated objects. As a result, false positives increase and AP score decreases. Incomplete Box Suppression (IBS) is proposed to reduce these kind of problems.

Essentially, IBS has the same principle with the NMS algorithm: finding the overlapped bounding boxes, selecting the box with the highest confidence value, and suppressing the others. While NMS uses a simple Intersection over Union(IoU) threshold to find overlaps, in the IBS, overlapping focal regions and object bounding boxes are both utilized to decide which box to suppress.

Let $C_i$ and $B_{ij}$ be the $i^{th}$ focal region coordinates and the $j^{th}$ box coordinates in that region.

- First step is to calculating IoU between focal region $C_i$ and other focal regions to find $C_i$'s overlaps. Overlapping focal regions are obtained after applying a threshold to calculated IoUs.
- Second step is to clip objects box coordinates in overlapping focal regions to the $i^{th}$ focal region's coordinates and gather the boxes whose areas are greater than zero.
- Final step is to calculate IoU between clipped boxes and the $B_{ij}$. If any of the IoU scores are greater than the selected threshold, the $B_{ij}$ is suppressed.

The IoU threshold for focal regions is experimentally selected as $0.05$, and IoU threshold for bounding boxes is, again, experimentally selected as $0.5$.
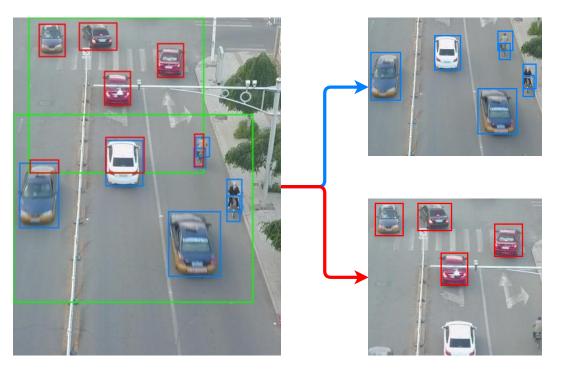
7

Figure 5: Incomplete Box Suppression (IBS). Overlapping blue bounding boxes suppress the incomplete red ones and decreases the false positives.

### 3.4.2 Non-max Suppression

Non-max suppression is applied to suppress overlapped detections after the merging of focal regions. Some of the overlapped focal regions contain same objects which causes duplicate box predictions. To mitigate this behavior, boxes with the highest confidences are selected and other boxes are suppressed. Intersection over union threshold for NMS is selected as $0.5$.

## 4 Experimental Results

### 4.1 Implementation Details

We implement Focus&Detect based on the publicly available MMDetection [55] and Pytorch. Generalized Focal Loss with Feature Pyramid Network is selected as the base detectors for both focus and detection stages. ResNet-50 and ResNeXt-101 are utilized as feature extraction networks in Focus stage and Detection stage, respectively. Focal detections are merged using NMS and IBS to obtain final predictions.

**Training phase** The input size of Focus stage is randomly sampled from $400 \times 1400$ to $1200 \times 1400$ in each step and samples are uniformly distributed for VisDrone dataset [13]. For UAVDT dataset [56], the input sizes of Focus stage and Detect stage are randomly sampled from $400 \times 1000$ to $800 \times 1000$ and $400 \times 800$ to $800 \times 800$, respectively. Flip augmentation is used with probability of $0.5$. For both Focus and Detection models, gradient descent with momentum, weight decay and learning rate scheduling is used. Both models are trained for 24 epochs. We set the initial learning rate to 0.01, at the 16th and 22th epoch, learning rate decreased to 0.001 and 0.0001. Beta parameter of momentum is selected as 0.9 for both models. Weight decay is applied with 0.0001 ratio. Both model leverages Synchronized Batch Normalization and Group Normalization on the backbone and FPN, respectively.

**Testing phase** In the experiments with VisDrone dataset, the input sizes of Focusing and Detection models are selected as $1200 \times 1400$ and $600 \times 1000$. On the other hand, in the experiments with UAVDT dataset the input sizes of Focusing and Detection models are selected as $600 \times 1000$ and $600 \times 800$. While merging detections of focal regions, NMS and IBS are applied. IoU threshold for NMS is $0.5$. After NMS, IBS is applied to reduce false positives caused by truncated objects in focal regions. IoU thresholds for IBS are selected as $0.05$ and $0.5$, where first threshold is for

overlapping focal regions and second threshold is for overlapping truncated objects in focal regions after clipping is applied.

## 4.2 Dataset and Evaluation Metric

In this work, we utilize VisDrone2021 Detection dataset [13] which consists of 6,471 images for training, 548 images for validation and 3,190 images for testing. It is an aerial image dataset which is obtained by drone-mounted cameras taken from 14 cities in China. Moreover, this dataset has 10 classes of objects which are non-uniformly distributed.

In order to assess the performance of our method, we use the evaluation protocol in MS COCO [33]. To be precise, we report AP, $AP_{50}$ and $AP_{75}$ scores, where AP shows the average precision for ten IoU thresholds whose range is from 0.5 to 0.95 and $AP_{50}$ is the average precision with the IoU threshold of 0.5. Similarly, $AP_{75}$ is the average precision with the IoU threshold of 0.75. In addition to these, we present $AP_S$, $AP_M$ and $AP_L$ scores, where $AP_S$, $AP_M$ and $AP_L$ represent the AP for objects with the area less than $32 \times 32$, less than $96 \times 96$ and larger than $96 \times 96$ pixels, respectively.

We also make experiments on UAVDT dataset [56], which is an aerial image dataset consists of around 41k frames with 840k bounding boxes. It has car, truck, and bus categories. The class distribution of the UAVDT dataset is extremely imbalanced where the truck and bus classes cover less than 5% of bounding boxes. Following the dataset authors' convention, we combine the three classes into one vehicle class and report PASCAL VOC AP score with an IoU threshold of 0.7 based on [56].

## 4.3 Results

We report results on VisDrone test set in Table 1, comparing our model to the baseline model. Focus&Detect significantly improves performance on pedestrian, person, bicycle classes which are mostly consist of small objects. On the other hand, baseline model yields a better performance in mAP@50 score on tricycle and awning-tricycle classes.

To have a fair comparison with the reported performance results of other state-of-the-art techniques, we also present our results on VisDrone validation set in Table 2. We compare our method with the other state-of-the-art region search based methods for object detection, in terms of number of images which are forwarded to detectors, evaluation metrics for precision and inference time per image.

The stated inference time results of our method and CRENet [16] are obtained on RTX 2080 Ti, while others are obtained on GTX 1080 Ti GPUs. We report the average inference time per image, since number of focal regions per image differs, similar to SAIC-FPN [12].

Table 1: Comparison between base model GFL and Focus&Detect, in terms of class-wise AP@50[%] scores on VisDrone **test-dev** dataset. Results indicate that the proposed Focus&Detect method significantly improves AP scores of the classes with the largest number of small objects, namely, "Pedestrian", "Person", and "Bicycle". In the experiments, GFL is tested in $2160 \times 3840$ resolution with flip augmentation, Focus&Detect is tested in $600 \times 1000$ resolution without any test-time augmentation. The highest AP scores for each column are printed in boldface.

| Method | All classes | Pedestrian | Person | Bicycle | Car | Van | Truck | Tricycle | Awning-tricycle | Bus | Motor |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GFL | 52.0 | 52.8 | 34.1 | 30.3 | **86.4** | 56.1 | 60.2 | **39.1** | **34.9** | 69.6 | **56.3** |
| Focus&Detect | **52.6** | **56.5** | **38.6** | **32.1** | 86.1 | **56.3** | **60.8** | 36.8 | 32.6 | **70.2** | 56.1 |

Results show that evenly image partition (EIP) method obtains the worst results [15], despite the high number of regions generated where it can be defined as the most straightforward approach among the region search based methods. Moreover, our method obtains 42.06 AP score on VisDrone validation set which outperforms other region search based object detection methods, reported in the literature. Although we set lower image resolution values as compared with CRENet, our method clearly surpasses CRENet in terms of $AP_S$ which happens to be the reported best result for $AP_S$.

Furthermore, we compare our method with the state-of-the-art methods based on various techniques, on VisDrone validation set as presented in Table 3. Results indicate that our method achieves 42.06 $AP$ and 44.64 $AP_{75}$ scores which are the best scores among state-of-the-art methods, reported in the literature.

We report results on UAVDT dataset in Table 4, comparing our model to the state-of-the-art reported methods on UAVDT dataset [56]. Compared to other methods Focus&Detect significantly improves performance. Results show that our framework outperforms FPN [20] and RetinaNet [26]. UAVDT dataset [56] does not solely focus on small

Table 2: Results on VisDrone validation set. Our method is compared with the other region search based object detection methods. "#img" is the number of images that the detectors are fed by, whereas the last column shows the inference time per image in seconds. We report the average inference time per image for our method, since number of focal regions per image differs. "SS" shows single-scale inference. * represents the result of EIP with the detector of Faster R-CNN and FPN, obtained by [15] which means partition image equally into six non-overlapping pieces to obtain regions to be focused. Bold values represent the column-wise best scores.

| Method | Backbone | Image Resolution | #Img | $AP[\%]$ | $AP_{50}[\%]$ | $AP_{75}[\%]$ | $AP_S[\%]$ | $AP_M[\%]$ | $AP_L[\%]$ | s / img |
|---|---|---|---|---|---|---|---|---|---|---|
| EIP* | ResNeXt-101 | $600 \times 1000$ | 3288 | 24.4 | 47.8 | 21.8 | 17.8 | 34.8 | 34.3 | 0.936 |
| ClusDet (SS) [15] | ResNeXt-101 | $600 \times 1000$ | 2716 | 28.4 | 53.2 | 26.4 | 19.1 | 40.8 | 54.4 | 0.773 |
| DMNet [47] | ResNeXt-101 | $600 \times 1000$ | 2736 | 29.4 | 49.3 | 30.6 | 21.6 | 41.0 | 56.9 | - |
| CRENet [16] | Hourglass-104 | $1024 \times 1024$ | 2337 | 33.7 | 54.3 | 33.5 | 25.6 | 45.3 | **58.7** | 0.901 |
| SAIC-FPN [12] | ResNeXt-101 | - | - | 35.7 | 63.0 | 35.1 | - | - | - | $0.252 \sim 2.568$ |
| AdaZoom [57] | ResNeXt-101 | - | - | 40.3 | **66.9** | 41.8 | - | - | - | - |
| **Focus&Detect** | ResNeXt-101 | $600 \times 1000$ | 9004 | **42.0** | 66.1 | **44.6** | **32.0** | **47.9** | 54.5 | 1.362 |

Table 3: Comparison of our method with other state-of-the-art methods for object detection on VisDrone validation set, where we report the results on original papers for other methods. Bold values represent the column-wise best scores.

| Method | Backbone | $AP[\%]$ | $AP_{50}[\%]$ | $AP_{75}[\%]$ |
|---|---|---|---|---|
| RRNet [45] | Hourglass | 32.92 | - | 31.33 |
| CRENet [16] | Hourglass-104 | 33.70 | 54.30 | 33.50 |
| DMNet [47] | ResNet-50 | 28.20 | 47.60 | 28.90 |
| CascadeNet [58] | ResNet-50 | 30.12 | 58.02 | 27.53 |
| GLSAN [59] | ResNet-50 | 30.70 | 55.40 | 30.00 |
| SAMFR [11] | ResNet-50 | 33.72 | 58.62 | 33.88 |
| AdaZoom (w/Faster R-CNN)[57] | ResNet-50 | 36.19 | 63.50 | 36.11 |
| MPFPN [60] | ResNet-101 | 29.05 | 54.38 | 26.99 |
| GLSAN [59] | ResNet-101 | 30.70 | 55.60 | 29.90 |
| ClusDet [15] | ResNeXt-101 | 32.40 | 56.20 | 31.60 |
| QueryDet [61] | ResNeXt-101 | 33.91 | 56.12 | 34.85 |
| SAIC-FPN [12] | ResNeXt-101 | 35.69 | 62.97 | 35.08 |
| AdaZoom (w/Faster R-CNN) [57] | ResNeXt-101 | 37.58 | 66.25 | 37.34 |
| AdaZoom (w/Cascade R-CNN) [57] | ResNeXt-101 | 40.33 | **66.94** | 41.77 |
| DREN [51] | ResNeXt-152 | 30.30 | - | - |
| **Focus&Detect** | ResNeXt-101 | **42.06** | 66.12 | **44.64** |

Table 4: Comparison of our method with other state-of-the-art methods for object detection on UAVDT test set, where we report the results on original papers for other methods.

| Method | $AP[\%]$ |
|---|---|
| RetinaNet [26] | 33.95 |
| LRF-Net [62] | 37.81 |
| FPN [20] | 49.05 |
| NDFT [63] | 52.03 |
| Focus&Detect | **54.16** |

objects. Contrary to VisDrone dataset, there are no annotation for very small objects and they are ignored in test time. Nevertheless, Focus&Detect surpasses the methods that are reported in the literature.
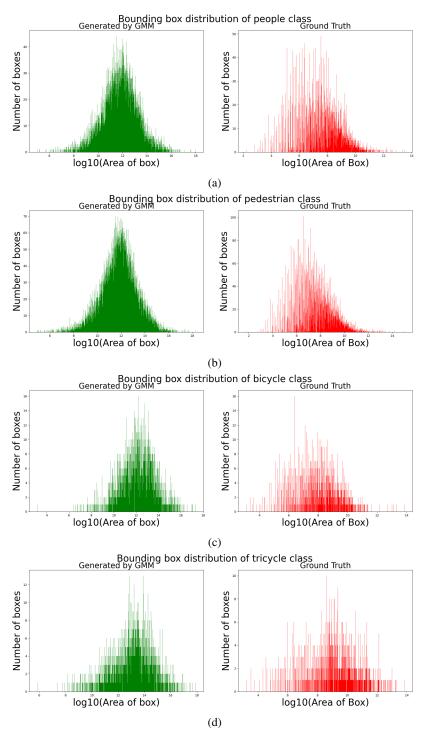
Figure 6: Effect of GMM clustering: We compare the areas of ground truth boxes and ground truth box areas of clusters obtained from GMM. We can clearly see that mean area is increased for each class. Moreover, distribution of boxes became normally distributed which clearly improve detection performance of models. (a,b,c,d) Distribution of bounding boxes for each class leverages from GMM clustering. Despite sizes of object clusters differ, bounding box areas are normalized.

Table 5: Comparison between base model GFL, Focus&Detect with IBS and Focus&Detect without IBS, in terms of class-wise AP@50[%] scores. GFL is tested in $2160 \times 3840$ resolution, Focus&Detect is tested in $600 \times 1000$ resolution.

| Method | All classes | Pedestrian | Person | Bicycle | Car | Van | Truck | Tricycle | Awning-tricycle | Bus | Motor |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GFL | 62.8 | 74.4 | 61.4 | 48.5 | 90.6 | 64.1 | 57.8 | 52.7 | 29.9 | 77.6 | 71.4 |
| F&D w/o IBS | 63.8 | 75.5 | 61.7 | 49.7 | 90.4 | 65.2 | 58.9 | 55.5 | 30.3 | 80.5 | 69.9 |
| F&D w/ IBS | 66.1 | 78.6 | 67.6 | 53.1 | 91.7 | 67.2 | 60.7 | 57.5 | 32.0 | 81.8 | 73.2 |

### 4.4 Ablation Study

We conduct ablation experiments in order to validate the contributions of GMM and IBS to the overall performance of object detection. To clarify the contributions, we report the class-wise results on VisDrone validation set in Table 5 and show the effect of GMM clustering in Figure 6.

**Effect of GMM Clustering**  Aerial images contain different sized objects depending on angle and altitude of the drone. Training an object detection model for aerial images is challenging as the data distribution does not contain all the different scales of objects. If ground truth bounding box areas are well distributed over real data distribution, object detector yields higher detection performance.

As seen in Figure 6, clustering objects has a normalizing effect on areas of bounding boxes. GMM clustering increases the mean box area which helps to improve performance on small objects.

We can compare effect of GMM clustering with resizing image to high resolution and multi-scale training. Increasing resolution shifts the mean of objects size. However, it does not normalize the bounding box areas. On the other hand Multi-scale training has similar effects, as it shifts means of bounding box areas and normalizes the box areas.

**Effect of IBS**  Results point out that F&D without IBS method improves particularly detection scores of 'bus' and 'tricycle' classes compared to GFL method.

In addition to this, proposed IBS method advances the detection performance of all classes by 2.3%. Moreover, IBS provides a performance boost on all classes, especially on classes of small objects such as 'person', 'bicycle', 'motor' and 'pedestrian'.

## 5 Conclusion

A two stage framework is proposed to solve small object detection problem in aerial images. The proposed method is region search based where we utilize a Gaussian Mixture Model to generate focal regions for object detection. GMM method has a normalization effect on GT box sizes as cropping and resizing the image to a fixed resolution relatively forces objects to an average size for each class. We also propose the Incomplete Box Suppression (IBS) method to mitigate the truncated box problem that arise while merging the target regions.

Results show that the proposed IBS method improves the detection performances of all classes, especially of small object classes. GMM clustering normalizes the object scales across regions and increases overall performance. Furthermore, our method achieves the state-of-the-art performance on VisDrone validation set and UAVDT test set comparing to other small object detection methods, reported in the literature. Moreover, our method obtains the best $AP_S$ score among all other methods, which indicates the positive impact of the proposed framework on small object detection.

## 6 Funding

## CRediT authorship contribution statement

**Onur Can Koyun:** Conceptualization, Methodology, Software, Writing - original draft. **Reyhan Kevser Keser:** Conceptualization, Methodology, Analysis and interpretation of data, Writing - original draft. **İbrahim Batuhan Akkaya:** Supervision, Writing - review & editing. **Behçet Uğur Töreyin:** Supervision, Writing - review & editing

## References

[1] Kang Tong, Yiquan Wu, and Fei Zhou. Recent advances in small object detection based on deep learning: A review. *Image and Vision Computing*, 97:103910, 2020.

[2] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 886–893. Ieee, 2005.

[3] David G Lowe. Object recognition from local scale-invariant features. In *Proceedings of the Seventh IEEE International Conference on Computer Vision*, volume 2, pages 1150–1157. Ieee, 1999.

[4] Yang Liu, Peng Sun, Nickolas Wergeles, and Yi Shang. A survey and performance evaluation of deep learning methods for small object detection. *Expert Systems with Applications*, page 114602, 2021.

[5] 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.

[6] Sujata Butte, Aleksandar Vakanski, Kasia Duellman, Haotian Wang, and Amin Mirkouei. Potato crop stress identification in aerial images using deep learning-based object detection. *arXiv preprint arXiv:2106.07770*, 2021.

[7] Ahlem Walha, Ali Wali, and Adel M Alimi. Moving object detection system in aerial video surveillance. In *International Conference on Advanced Concepts for Intelligent Vision Systems*, pages 310–320. Springer, 2013.

[8] Fan Zhang, Bo Du, Liangpei Zhang, and Miaozhong Xu. Weakly supervised learning based on coupled convolutional neural networks for aircraft detection. *IEEE Transactions on Geoscience and Remote Sensing*, 54(9):5553–5563, 2016.

[9] Gong Cheng, Peicheng Zhou, and Junwei Han. Rifd-cnn: Rotation-invariant and fisher discriminative convolutional neural networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2884–2893, 2016.

[10] Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. DOTA: A large-scale dataset for object detection in aerial images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3974–3983, 2018.

[11] Haoran Wang, Zexin Wang, Meixia Jia, Aijin Li, Tuo Feng, Wenhua Zhang, and Licheng Jiao. Spatial attention for multi-scale feature refinement for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[12] Jingkai Zhou, Chi-Man Vong, Qiong Liu, and Zhenyu Wang. Scale adaptive image cropping for UAV object detection. *Neurocomputing*, 366:305–313, 2019.

[13] Pengfei Zhu, Longyin Wen, Xiao Bian, Haibin Ling, and Qinghua Hu. Vision meets drones: A challenge. *arXiv preprint arXiv:1804.07437*, 2018.

[14] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Haibin Ling, Qinghua Hu, Qinqin Nie, Hao Cheng, Chenfeng Liu, Xiaoyu Liu, et al. Visdrone-det2018: The vision meets drone object detection in image challenge results. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.

[15] Fan Yang, Heng Fan, Peng Chu, Erik Blasch, and Haibin Ling. Clustered object detection in aerial images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8311–8320, 2019.

[16] Yi Wang, Youlong Yang, and Xi Zhao. Object detection using clustering algorithm adaptive searching regions in aerial images. In *European Conference on Computer Vision*, pages 651–664. Springer, 2020.

[17] 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.

[18] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *European Conference on Computer Vision*, pages 483–499. Springer, 2016.

[19] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1492–1500, 2017.

[20] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017.

[21] 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.

[22] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018.

[23] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.

[24] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016.

[25] Xiang Li, Wenhai Wang, Lijun Wu, Shuo Chen, Xiaolin Hu, Jun Li, Jinhui Tang, and Jian Yang. Generalized focal loss: Learning qualified and distributed bounding boxes for dense object detection. *arXiv preprint arXiv:2006.04388*, 2020.

[26] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[27] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9627–9636, 2019.

[28] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6569–6578, 2019.

[29] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S Davis. Soft-NMS–improving object detection with one line of code. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5561–5569, 2017.

[30] Alexander Neubeck and Luc Van Gool. Efficient non-maximum suppression. In *18th International Conference on Pattern Recognition (ICPR'06)*, volume 3, pages 850–855. IEEE, 2006.

[31] Xiyang Dai, Yinpeng Chen, Bin Xiao, Dongdong Chen, Mengchen Liu, Lu Yuan, and Lei Zhang. Dynamic head: Unifying object detection heads with attentions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7373–7382, 2021.

[32] Tingting Liang, Xiaojie Chu, Yudong Liu, Yongtao Wang, Zhi Tang, Wei Chu, Jingdong Chen, and Haibing Ling. CBNetV2: A composite backbone network architecture for object detection. *arXiv preprint arXiv:2107.00420*, 2021.

[33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014.

[34] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *European Conference on Computer Vision*, pages 213–229. Springer, 2020.

[35] Junhyug Noh, Wonho Bae, Wonhee Lee, Jinhwan Seo, and Gunhee Kim. Better to follow, follow to be better: Towards precise supervision of feature super-resolution for small object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9725–9734, 2019.

[36] Jakaria Rabbi, Nilanjan Ray, Matthias Schubert, Subir Chowdhury, and Dennis Chao. Small-object detection in remote sensing images with end-to-end edge-enhanced gan and object detector network. *Remote Sensing*, 12(9):1432, 2020.

[37] Heming Zhang, Xiaolong Wang, Jingwen Zhu, and C-C Jay Kuo. Fast face detection on mobile devices by leveraging global and local facial characteristics. *Signal Processing: Image Communication*, 78:1–8, 2019.

[38] Seyed Mojtaba Marvasti-Zadeh, Javad Khaghani, Hossein Ghanei-Yakhdan, Shohreh Kasaei, and Li Cheng. COMET: context-aware IoU-guided network for small object tracking. In *Proceedings of the Asian Conference on Computer Vision*, 2020.

[39] Chang Sun, Yibo Ai, Sheng Wang, and Weidong Zhang. Mask-guided SSD for small-object detection. *Applied Intelligence*, 51(6):3311–3322, 2021.

[40] Haodong Pan, Jue Jiang, and Guangfeng Chen. TDFSSD: Top-down feature fusion single shot multibox detector. *Signal Processing: Image Communication*, 89:115987, 2020.

[41] Qunjie Yin, Wenzhu Yang, Mengying Ran, and Sile Wang. FD-SSD: An improved SSD object detection algorithm based on feature fusion and dilated convolution. *Signal Processing: Image Communication*, page 116402, 2021.

[42] Christian Eggert, Dan Zecha, Stephan Brehm, and Rainer Lienhart. Improving small object proposals for company logo detection. In *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, pages 167–174, 2017.

[43] Mate Kisantal, Zbigniew Wojna, Jakub Murawski, Jacek Naruniec, and Kyunghyun Cho. Augmentation for small object detection. *arXiv preprint arXiv:1902.07296*, 2019.

[44] Xuerui Dai. HybridNet: A fast vehicle detection system for autonomous driving. *Signal Processing: Image Communication*, 70:79–88, 2019.

[45] Changrui Chen, Yu Zhang, Qingxuan Lv, Shuo Wei, Xiaorui Wang, Xin Sun, and Junyu Dong. Rrnet: A hybrid detector for object detection in drone-captured images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[46] Sungeun Hong, Sungil Kang, and Donghyeon Cho. Patch-level augmentation for object detection in aerial images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[47] Changlin Li, Taojiannan Yang, Sijie Zhu, Chen Chen, and Shanyue Guan. Density map guided object detection in aerial images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 190–191, 2020.

[48] F Ozge Unel, Burak O Ozkalayci, and Cevahir Cigla. The power of tiling for small object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.

[49] Ziyang Tang, Xiang Liu, Guangyu Shen, and Baijian Yang. Penet: object detection using points estimation in aerial images. *arXiv preprint arXiv:2001.08247*, 2020.

[50] Zhiwei Wei, Chenzhen Duan, Xinghao Song, Ye Tian, and Hongpeng Wang. AMRNet: Chips augmentation in aerial images object detection. *arXiv preprint arXiv:2009.07168*, 2020.

[51] Junyi Zhang, Junying Huang, Xuankun Chen, and Dongyu Zhang. How to fully exploit the abilities of aerial image detectors. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[52] Xizhou Zhu, Han Hu, Stephen Lin, and Jifeng Dai. Deformable convnets v2: More deformable, better results. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9308–9316, 2019.

[53] Jifeng Dai, Haozhi Qi, Yuwen Xiong, Yi Li, Guodong Zhang, Han Hu, and Yichen Wei. Deformable convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 764–773, 2017.

[54] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010.

[55] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open MMLab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.

[56] Dawei Du, Yuankai Qi, Hongyang Yu, Yifan Yang, Kaiwen Duan, Guorong Li, Weigang Zhang, Qingming Huang, and Qi Tian. The unmanned aerial vehicle benchmark: Object detection and tracking. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 370–386, 2018.

[57] Jingtao Xu, Yali Li, and Shengjin Wang. AdaZoom: Adaptive zoom network for multi-scale object detection in large scenes. *arXiv preprint arXiv:2106.10409*, 2021.

[58] Xindi Zhang, Ebroul Izquierdo, and Krishna Chandramouli. Dense and small object detection in UAV vision based on cascade network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[59] Sutao Deng, Shuai Li, Ke Xie, Wenfeng Song, Xiao Liao, Aimin Hao, and Hong Qin. A global-local self-adaptive network for drone-view object detection. *IEEE Transactions on Image Processing*, 30:1556–1569, 2020.

[60] Yingjie Liu, Fengbao Yang, and Peng Hu. Small-object detection in UAV-captured images via multi-branch parallel feature pyramid networks. *IEEE Access*, 8:145740–145750, 2020.

[61] Chenhongyi Yang, Zehao Huang, and Naiyan Wang. QueryDet: Cascaded sparse query for accelerating high-resolution small object detection. *arXiv preprint arXiv:2103.09136*, 2021.

[62] Tiancai Wang, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. Learning rich features at high-speed for single-shot object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1971–1980, 2019.

[63] Zhenyu Wu, Karthik Suresh, Priya Narayanan, Hongyu Xu, Heesung Kwon, and Zhangyang Wang. Delving into robust object detection from unmanned aerial vehicles: A deep nuisance disentanglement approach. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1201–1210, 2019.