

---

# PERCEIVING TRAFFIC FROM AERIAL IMAGES

---

**George Adaimi**

VITA, EPFL

Switzerland

george.adaimi@epfl.ch

**Sven Kreiss**

VITA, EPFL

Switzerland

sven.kreiss@epfl.ch

**Alexandre Alahi**

VITA, EPFL

Switzerland

alexandre.alahi@epfl.ch

## ABSTRACT

Drones or UAVs, equipped with different sensors, have been deployed in many places especially for urban traffic monitoring or last-mile delivery. It provides the ability to control the different aspects of traffic given real-time observations, an important pillar for the future of transportation and smart cities. With the increasing use of such machines, many previous state-of-the-art object detectors, who have achieved high performance on front facing cameras, are being used on UAV datasets. When applied to high-resolution aerial images captured from such datasets, they fail to generalize to the wide range of objects' scales. In order to address this limitation, we propose an object detection method called Butterfly Detector that is tailored to detect objects in aerial images. We extend the concept of fields and introduce butterfly fields, a type of composite field that describes the spatial information of output features as well as the scale of the detected object. To overcome occlusion and viewing angle variations that can hinder the localization process, we employ a voting mechanism between related butterfly vectors pointing to the object center. We evaluate our Butterfly Detector on two publicly available UAV datasets (UAVDT and VisDrone2019) and show that it outperforms previous state-of-the-art methods while remaining real-time.

**Keywords** Traffic Monitoring · UAV · Object Detection · Aerial Images

## 1 Introduction

Smart cities have been increasing around the world especially with the increasing importance of sustainable development and the effect of AI on our everyday life. It is projected also that the market for such cities will keep on growing [1]. For a successful implementation of smart cities, a large amount of data is required to better understand its different aspects such as congestion and other transportation network conditions. Visual information have long been collected through the use of many static cameras placed at different locations. Installing such cameras with the required infrastructure is costly and fails to monitor a very large area. Moreover, certain unexpected scenarios require visual information at un-monitored location, such as in the case of huge congestions caused

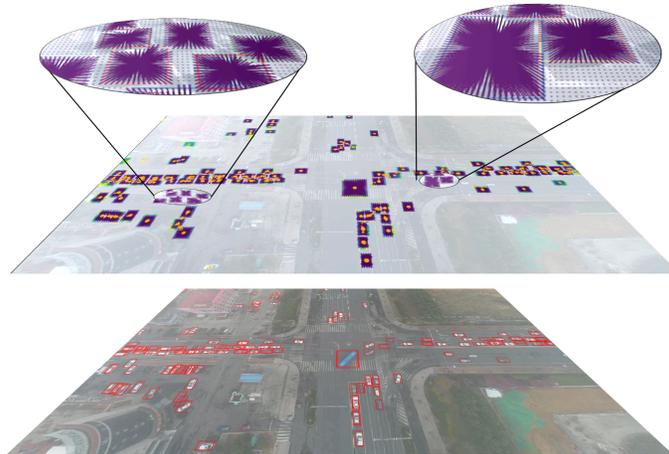


Figure 1: Butterfly fields outputted by our method overlaid on top of an aerial image (bottom). The purple vector components are pointing towards the center (orange circle). Using the different components of our butterfly fields allows us to better detect and localize the vehicles in the image

by an accident. As a result, data collected from unmanned aerial vehicles (UAV) or drones have been increasing in popularity for urban traffic monitoring. One such example is the introduction of a new drone-collected large-scale dataset called pNEUMA [2] which aims to help researchers better study and model traffic congestion. Collecting such large datasets requires analyzing large amounts of objects across many images. This work was traditionally outsourced to manual labellers such as Amazon Mechanical Turk (MTurk) [3]. With the huge success of deep learning, especially in the field of computer vision, most of the labelling work has become automated. In addition to using UAVs for dataset collection, many people believe that drones is the solution to increase the efficiency of last-mile deliveries [4, 5, 6, 7, 8]. This requires the drone to be able to visually understand its surrounding using different sensors. Thus, paradigms for different vision tasks are still challenged with an arms race of methods every year. One of the most prominent vision tasks is *object detection*. From traditional methods such as sliding window and deformable parts model (DPM) [9] to deep neural networks with different components [10], the object detection race is still on.

Although general object detectors have achieved promising results on front facing cameras [11, 12], they fail to perform as well when used from a drone/UAV point-of-view. Aerial images capture a large scene made of hundreds of objects occupying a tiny portion of the image with very limited number of pixels (Figure 1). Hence, this adds an exciting challenge to the problem of object detection. In this work, we propose a method tailored to detect objects in aerial images. Our method targets applications such as traffic analysis, smart cities, and more broadly digital twins to name a few.

With the prevalence of UAVs with high resolution cameras, two new datasets were recently released called UAVDT [13] and VisDrone2019 [14]. A common practice is to re-train and adapt some of the design choices of existing object detectors (*e.g.*, Yolo [15]) to the new distribution of the data. For instance, Wu et al. [16] adapted the sizes of the anchor boxes of Faster R-CNN [17] to improve its results on aerial images. Ding et al. [18] also introduced an ROI transformer to replace

the conventional ROI and deal with a common challenge in UAV images, variations in objects' orientations. Yet, state-of-the-art methods still suffer from detecting tiny objects of any rotation.

Recent object detectors, using deep neural networks, can be categorized as anchor-based approaches and anchor-free approaches. On one hand, anchor-based methods leverage strong priors on scale and aspect ratio by pre-defining specific bounding boxes as initial detections. Rather than directly detecting the bounding boxes, the objects are located relative to pre-defined boxes obtained by analyzing the training dataset. Hence, they do not generalize well to objects of different scales or to images of different resolutions, especially when used with different datasets. On the other hand, anchor-free methods focus on locating corners and centers of objects directly with the use of prior knowledge. They are very challenging to train due to the sparsity of their output feature map. This is especially true in the case of objects in high-resolution images, a common characteristic of aerial images.

To overcome the limitations of both anchor-based and anchor-free approaches, we introduce a novel detector, called Butterfly Detector (BD), that extends and adapts field theory, which has been gaining interest in the 2D human pose estimation task [19, 20, 21]. In previous anchor-free models, only a specific feature in the output feature map of a model is responsible to localize the center or corner of an object. In our work, we make use of fields, which we call *butterfly fields*, to localize the center and predict the width and height of an object. Each output feature that contains information about the object has a set of field components spatially pointing towards the center. Rather than implementing a naive approach of using all predicted bounding boxes at every output feature, we implement a voting mechanism leveraging the direction and information of related points in a field to improve localization as well as obtain a significant reduction in false positives. We also improved our precision, especially for small objects, by relying on a sub-pixel localization. Our experiments show that we outperform previous state-of-the-art methods on UAVDT and VisDrone2019 dataset. We especially achieve a boost in performance on images that are recorded at high altitudes (Table 1) showcasing the importance of our method on aerial images. The code will be made available<sup>1</sup>. Our method can also be tested using the video-labelling tool, UltimateLabeling<sup>2</sup>, where it is used to detect vehicles.

## 2 Related Work

Object detection can be done at 2D scale [17, 22, 15, 23] or 3D scale [24, 25, 26, 27]. Since we are dealing with aerial images, we focus on 2D detection and further divide these methods into traditional detectors, anchor-based and anchor-free deep learning detectors.

### 2.1 Traditional detectors for aerial images

During the rise of satellites and aerial imagery, a lot of effort was put in developing methods to detect different objects [28]. This was quite challenging because of many different variation such as the small object size and different weather conditions. One of the earliest form of object detection relied on using a hand-crafted template per class and measuring its similarity at every location in the image using the sliding window technique. A lot of advancements have been done to improve

---

<sup>1</sup><https://github.com/vita-epfl/butterflydetector>

<sup>2</sup><https://github.com/alexandre01/UltimateLabeling>

the templates used, from rigid templates [29, 30] to deformable ones [31, 32]. However, these methods are difficult to generalize to different shapes and orientations. To solve this, researchers moved toward extracting different features that help in localization and classification, such as color intensities and color gradients. One of the most interesting work in detection was done by [33] where they introduced histogram of oriented gradients (HoG) for human detection. The features extracted from the image are the distribution of oriented gradients, which allows the detection of edges. Many work since then extended this idea and further improved it to deal with different variations that occur while detecting objects. For instance, rotation-invariant HoG features [34] were developed to make these features resilient to rotational changes while deformable HoG [35] were introduced to be resilient against occlusion and shape changes. Although pioneering works relied on handcrafted features and often dealt with satellite images [36, 37, 38, 39, 40, 41], we will focus on reviewing recent deep learning based methods tailored for UAV images.

## 2.2 Anchor-based detector for UAV images

When AlexNet [42] was first used for vehicle detection in aerial images, it was able to surpass all traditional methods. Since then, deep convolutional neural networks have dominated object detection with the release of Faster R-CNN [17], Mask R-CNN [22], YOLO [15], and SSD [23]. Although these methods achieved good performance on natural images (MS COCO [43] and PASCAL VOC [44]), they perform poorly on high-resolution aerial images. Since most of these methods rely on anchor boxes and objects that are rather big, they require a lot of modifications to improve their performance on a different task such as UAV object detection.

Rather than re-inventing the wheel, many works have extended Faster R-CNN [17], R-FCN [45], SSD [23], and YOLO [15] for vehicle detection [46, 47, 48, 49, 50, 51] by proposing many adaptations to overcome challenges characteristic of aerial images. These adaptations include changing the number and size of anchor boxes, increasing the resolution, as well as making use of skip connections in the network. Tang et al. [52] proposed replacing the RPN module of Faster R-CNN with a hyper region proposal network (HRPN) that makes use of hierarchical feature maps. Deng et al. [53] extended this work by learning vehicle attributes as well as location and type.

Wu et al. [16] proposed augmenting an off-the-shelf detector to also learn different attributes for every input such as altitude, weather, and viewing angle. This is done in an adversarial manner to disentangle these variations from the extracted features. Another work by Cai et al. [54] uses an anchor-free architecture with background and foreground attention modules as well as multi-level features to improve object detection. In order to deal with small objects, Duan et al. [55] introduced a channel-aware de-convolution layer to exploit features from multiple channels of different layers. They also developed a Multi-RPN module that performs multiple detection at different layers simultaneously. In contrast, Li et al. [56] combined bottom-up and top-down attention mechanisms to extract more discriminative features. By also counting the number of objects in a scene, they were able to show that these two tasks aid each other and lead to a boost in performance. Other methods utilized clustering [57], feature fusion network [58], or rotation-invariant cascaded forest [59] to perform vehicle detection. All these methods, make use of either extra annotations, attention, or the fusion of multiple features from different layers to deal with different scales. This makes both the training and inference more complicated and slow.

### 2.3 Anchor-free Object Detectors

Due to the drawbacks of anchor-based object detectors, there has been recent work focusing on developing anchor-free detectors. Such methods detect the corners and centers of bounding boxes by outputting a heatmap over the image highlighting their location [60, 61, 11, 62, 63, 12]. Point Linking Network (PLN) [61] detects the four corners and centers of bounding boxes. Then each corner predicts which output feature cell has a high probability of being the center. Then a bounding box is formed from each pair of center and corner. PLN is different than our Butterfly detector since we do not only make use of corners but rather every point inside the bounding box. Moreover, instead of predicting a probability over all feature cells, each point predicts a vector pointing to the center, allowing us to have a floating point precision of the center location. Law et al. [11, 62], on the contrary, predict only two opposite corners and use associative embedding to group them together. Since corners usually lie outside an object, they might not contain a lot of appearance features. In order to remedy this problem, Zhou et al. [63] suggested using points that lie at the extreme edges of an object. They were able to show promising results, but this limits the use of this method on datasets that contain object masks. CenterNet [12] showed performance improvement by detecting the center as well as the size of the object. This method also shares many similarities with DenseBox [64], GuidedAnchoring [65], FAFS [66], FCOS [67], and FoveaBox [68] where they output two types of maps: (i) a confidence (heatmap) per object class predicting the center of objects, and (ii) a shape map predicting the bounding boxes. Instead of directly predicting a heatmap of centers, we propose to utilize the field vectors to point to object centers. Objects in aerial images experience huge variations in orientation compared to general object datasets (MS-COCO [43]). We solve this problem by having many output feature cells point towards the center, decreasing the reliance on specific features such as only the center of an object. To take into account the consensus of all vectors, we develop a voting mechanism to locate the center of an object. This mechanism provides us with a two fold advantage: we not only suppress overlapping predictions, but also accumulate the information thereby providing a more refined output.

Instead of directly predicting a heatmap of centers, we use vectors from different output cells to point to object centers. A similar idea is used for instance segmentation [69, 70, 71, 72] where the vectors predicted from different output cells is used to cluster related pixels. In contrast, rather than using these vectors to group feature cells, we make explicit use of where these vectors are pointing to localize the center of an object. We thus develop a voting mechanism to refine our predictions by allowing different vectors to agree to a final prediction, allowing us to deal with occlusion as well as scale variations.

## 3 Method

The goal of our method is to detect objects (such as cars, buses, and pedestrians) in aerial images. We propose an anchor-free method that overcomes the limitations of current anchor-based and anchor-free methods. Existing anchor-free detectors output a sparse heatmap wherein each cell is in charge of only locally detecting either a corner or a center of a bounding box. This is typically challenging for a network due to the fact that it needs to output a low confidence even for features that contain information about the object but do not lie at the specific corner of interest. Thus, we extend fields from 2D pose estimation to object detection and introduce **Butterfly Fields** which are named based on the shape they form when drawn over an image (Figure 1).

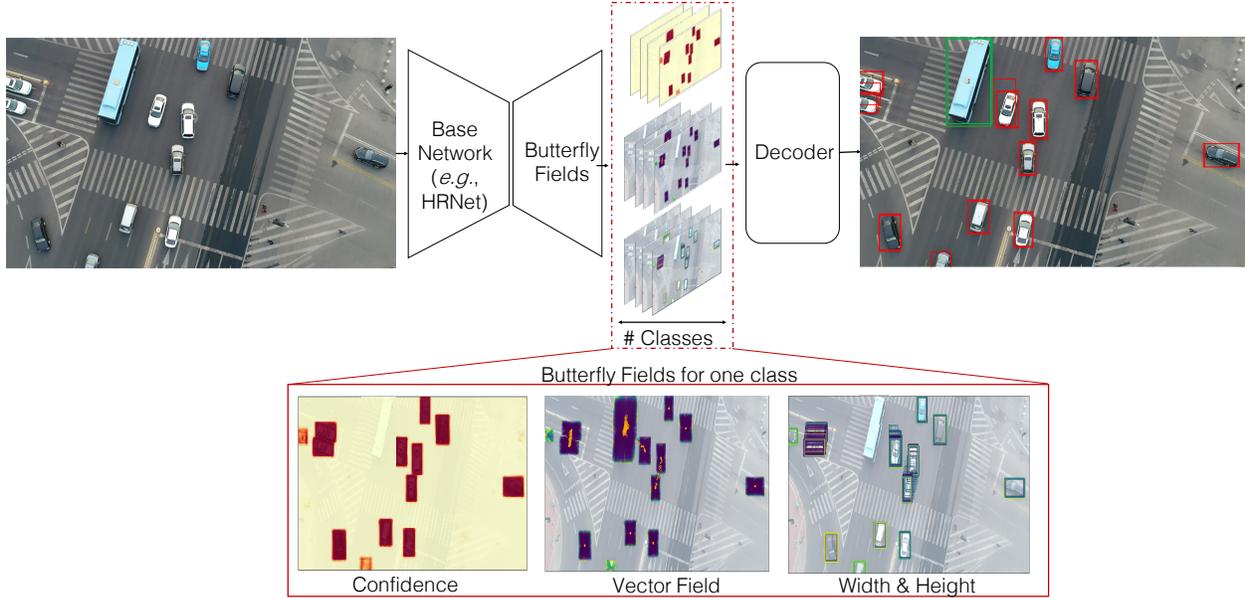


Figure 2: Model architecture showing the different components of our butterfly fields. The features extracted by the base-network are used by the head responsible to output the butterfly fields. The top component, represented by a heatmap over the input, indicates the confidence scalar component. The middle component is made up of vectors where each activated feature cell is pointing towards the center (orange circle). The bottom component shows the width and height predicted by each activated output cell. The resulting fields are then decoded to output the final bounding box predictions (right image)

Figure 2 illustrates the overall method. A base network, such as HRNet [73], is used to extract features from an aerial image. These features are then processed by the head responsible in outputting the composite fields, made up of a confidence, width and height of the bounding box, as well as vectors pointing towards the center of every object. In the following sections, we highlight the main parts of our method and the advantage of using fields for object detection compared to previous methods.

### 3.1 Butterfly Fields

General anchor-free object detectors usually assign a single feature cell to be responsible to detect the center of the object by predicting a confidence as well as the offset of the center relative to this cell [62, 61]. Some methods also detect other attributes such as the width and height [12]. This forces the neural network to output a sparse feature map with very few and specific activated location while reducing all cells around them. This task is extremely challenging in high-resolution images, where many cells may contain similar information. To overcome this drawback, we introduce **Butterfly Fields**. The goal of these fields is to localize the center of the object, width, and height from every feature cell that contains information about the object of interest. The difference between other anchor-free methods is explained in Figure 3. Predicting the center, height, and width, instead of corners does not require association between the corners, which would make inference slower.

For instance, if  $N$  top-left corners and  $N$  bottom-right corners were detected, this requires  $N^2$  computations to cluster the corners and form the bounding boxes.

Butterfly fields are composite fields made up of a scalar component for confidence, a vector component that points toward the closest center of an object, and two scalar components indicating the width and height of the predicted object. This can be represented over the feature map  $f$  as:

$$f_{ij}^c = \{p_{ij}^c, x_{ij}^c, y_{ij}^c, w_{ij}^c, h_{ij}^c\} \quad (1)$$

where  $i$  and  $j$  index the output feature cell, and  $c$  represents the object class.  $p$  is the confidence of the cell in its prediction,  $x$  and  $y$  form the vector from the center of the cell to the center of the object.  $w$  and  $h$  are the the width and height respectively of the referred object. In the case of overlapping objects of same class, the vector components will point to the closest center since the datasets used do not have any information about relative object positions.

As can be seen from Equation 1, a composite field is predicted for every object category. This allows the network to output different fields that are specialized for specific objects with specific aspect ratios. For instance, a model can learn that a pedestrian has a specific width to height ratio and thus the output feature map specific for this class would capture this characteristic. On the contrary, previous keypoint-based object detectors predict a common width, height, and even offset for all categories. Our modification allows us to obtain specialized predictors for different classes.

### 3.2 Loss functions

In order to train the model, a ground-truth feature map is built based on the bounding box annotations provided by the dataset. At every feature cell inside a bounding box, the butterfly field components are set such that the confidence  $p$  is 1.  $w$  and  $h$  are set to be the logarithm of the down-scaled width and height respectively in Eq 1. The vector component ( $x$  and  $y$ ) are set to point to the center of the object from the center of the corresponding feature cell. All other feature cells that do not contain any object have their confidence component  $p$  set to zero. We also compare the performance of our model when different number of cells around the center are chosen to predict the bounding box. For instance, we test on our model when using a 4x4 window around the center Figure 3b and when using all feature cells inside the bounding box Figure 3c.

Two main losses are used in order to train our network. The scalar components that represent the confidence are trained with independent binary cross-entropy, while the components of the field related to the width and height are trained using an L1 Loss. To train the vector components of the butterfly field, we make use of the Laplace Loss [74, 75]:

$$L = \frac{|x - \mu|}{b} + \log 2b \quad (2)$$

The model learns  $b$  without any supervision and is used to attenuate the gradients based on the model’s confidence in its predictions. It is important to note that we do not deal with the imbalance of negative and positive samples compared to other anchor-free object detectors that make use of a modified focal loss [76]. The focal loss was tested and did not show any improvement compared to using the normal binary cross-entropy.

### 3.3 Decoding: Consensus of every point in a field

During inference, given the fields over the output feature, we extract the set of bounding boxes with their confidences and appropriate classes. In order to deal with the huge variation of objects, we develop a decoding procedure that scales according to the height and width of every object.

For every class-specific feature output, we build a high resolution confidence map using both the confidence and vector components of the butterfly fields. This produces a heatmap over the input image where the intensity of a region corresponds to its probability of being the center of an object. We propose the following adaptive function to build our heatmap given the output feature map:

$$f(x, y) = \sum_{ij} \frac{p_{ij}}{\chi_{ij}} \exp \left( -\frac{1}{2} \left( \frac{x - x_{ij}}{\sigma_{ij}^x} \right)^2 - \frac{1}{2} \left( \frac{y - y_{ij}}{\sigma_{ij}^y} \right)^2 \right) \quad (3)$$

where  $x$  and  $y$  index a pixel in the high-resolution confidence map, and  $i$  and  $j$  vary over the output feature cells with confidences greater than 0.1. Each output cell predicts a center of an object  $(x_{ij}, y_{ij})$  and confidence  $p_{ij}$ . Instead of adding the confidences only at the predicted center, Equation 3 uses an unnormalized 2D Gaussian to spread this confidence in a neighborhood of size  $\sigma_{ij}$ . A variable number of output cells can predict the same object center and thus we need to deal with the variable normalization of the accumulated confidence. Since summing all the  $p_{ij}$  in Equation 3 will result in confidences greater than 1, we use  $\chi_{ij}$  to normalize the confidence based on the number of expected vectors voting for the same center. This value scales with the number of cells localizing an object center (related to the number of output cells inside the bounding box). In the case of using 16 cells around the center (4x4 window),  $\chi_{ij} = 16 \forall i, j$ . When using all the feature cells inside the bounding box to predict the center,  $\chi_{ij}$  is the area of the corresponding box. We define  $\sigma_{ij}^x$  and  $\sigma_{ij}^y$  as functions of the width and height respectively:

$$\sigma_{ij}^x = \max \left( 2, \frac{w_{ij}}{\rho} \right) \quad (4) \quad \sigma_{ij}^y = \max \left( 2, \frac{h_{ij}}{\rho} \right) \quad (5)$$

These equations can be interpreted as the maximum number of allowed error pixels for predicted bounding box center. This depends on the size of the object since the error in prediction for small objects is more costly than the prediction for big objects.  $\rho$  is a hyperparameter that defines the amount of error allowed relative to the size of the object. It can be a single number or a distinct number for each class. We set the minimum value for  $\sigma$  to be 2 pixels (shown in Equation 4 and 5).

Based on this voting mechanism, the resulting high-resolution confidence map contains peaks highlighting the probability of the existence of object centers. In addition, another consensus occurs when calculating both the width and height. They are calculated using a weighted average of all predictions following this equation:

$$w(x, y) = \frac{\sum_{ij} p_{ij} w_{ij}}{\sum_{ij} p_{ij}} \quad (6) \quad h(x, y) = \frac{\sum_{ij} p_{ij} h_{ij}}{\sum_{ij} p_{ij}} \quad (7)$$

where  $i$  and  $j$  are the indices of all feature cells that have their vector fields pointing towards a specific feature located at  $x$  and  $y$ .

Figure 3 illustrates the two key differences of our method with respect to previous works. Firstly, the confidence map of previous methods is usually dense only around the center of the object (Figure 3a). Such a sparse feature map results in difficulty in training. To tackle this, our method utilizes a greater number of feature cells within each bounding box. Each feature cell falling inside

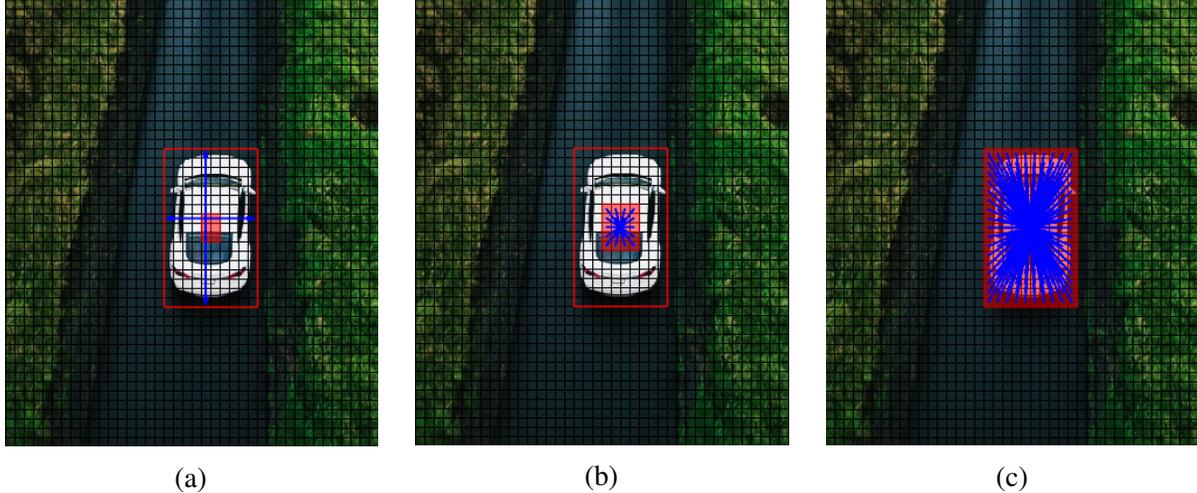


Figure 3: A visualization of the difference between anchor-free methods (a) [64, 65, 66, 67, 68] and our butterfly detector (b-c). (b) shows our method using a 4x4 window around to point towards the center. (c) shows our method using all grid cells inside the bounding box. Instead of having a map indicating the center, we make use of fields over the image where each vector is pointing towards the center to localize it. The voting mechanism aggregates information from multiple cells to obtain a refined prediction

the bounding box predicts a vector pointing toward the object center (Figure 3b-3c). Moreover, the feature cell confidence  $p$  indicates the confidence of the pointed vector rather than the probability of this cell being the center of an object. A map is also built where each feature cell predicts the height and width at the location pointed by the vector. Secondly, previous methods rely on Non-Maximum Suppression to filter their predictions. Instead, we propose our voting mechanism that allows the network to overcome many challenges such as occlusion. Consider the case in which a confidence map detecting the center is used, such as in CenterNet [12]; if more than half of the object of interest is occluded, the method would find difficulties detecting the center. Since we leverage cues from different feature cells to detect the center, we are able to overcome such shortcomings. For instance, by detecting only the trunk of a car, our method is able to correctly output a tight bounding box as shown in the qualitative results (Figure 5).

### 3.4 Decoding: Sub-pixel Localization

Upon building a high-resolution confidence map, cells with high confidence represent the center of objects. A loss of spatial precision occurs when using the center of cells as the centers. As can be seen in Figure 4a, using the center of the cell with the most Gaussian overlap might lead to an average between two predictions. This will lead to a loss of precision and will affect the detection of small objects. To remedy this issue, we make use of the floating-precision of our vector components. Instead of taking the center of the highest confident cell as the center of an object (red dot in Figure 4a), we take the exact location the vector is pointing to inside this cell (red dot in Figure 4b). This method leads to a boost in performance, especially for the VisDrone2019 since it contains many small objects, such as people and bicycles (see Section 4).

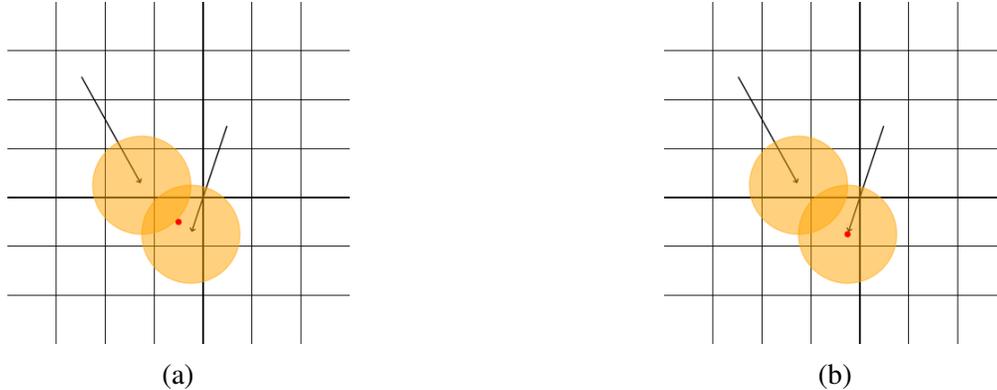


Figure 4: A visualization comparing using sub-pixel decoding (b) to using directly the high resolution map (a). The yellow circles represent the Gaussian used to fill the high-confidence resolution map (Eq. 3). The red circle is the resulting object center since it should be located in the cell with a high confidence

## 4 Experiments

To evaluate our method, we study the performance of the butterfly detector on two publicly available UAV datasets: UAVDT [13] and VisDrone2019 [14]. In order to showcase the importance of the different parts of our method, we also perform a detailed ablation study.

Our experiments show the advantage of using composite fields for object detection. Combined with a high resolution confidence accumulation procedure and a sub-pixel localization, we demonstrate state-of-the-art-performance for object detection in aerial images.

### 4.1 Datasets

**UAVDT** is a dataset collected using a UAV platform at different locations, altitudes, weather conditions, as well as viewing angle. This makes it a challenging dataset for both tracking and detection. The dataset also includes annotations for the previously listed variations in addition to the amount of occlusion. It contains three different classes: car, truck, bus. The detection aspect of this dataset is made up of around 41k frames from 50 distinct videos, 20 of which are used for testing.

**VisDrone 2019** is a new dataset that is also collected using a UAV platform. Compared to UAVDT, the detection aspect of this dataset is relatively small but contains much more categories. The 8.5k images include 10 classes: pedestrian, people, bicycle, car, van, truck, tricycle, awning-tricycle, bus, and motor. The small number of images and big number of categories introduces another challenge for our method. Moreover, the difference between pedestrian and people category is that pedestrians are people who are walking or standing still.

### 4.2 Implementation and Training Details

Our model uses an ImageNet pre-trained HRNet as a base-network since it recently showed good performance for keypoint estimation, a task that shares some similarities to our method. We

augment the base-network with a head composed of  $1 \times 1$  convolutional layers responsible to output the different components of our butterfly fields.

During training, we apply different standard data augmentation techniques. After normalizing the images, we randomly flip the input horizontally. The resulting image is then rescaled by a scale between 0.4 and 2.0. We subsequently extract a random  $512 \times 512$  crop that also undergoes a random  $90^\circ$  rotation.

As mentioned previously, the confidence map is trained using a binary cross entropy loss. The vector components of our fields are trained using a Laplace loss while all other scalar components, such as width and height, are trained using an L1 loss. In order to find the best prefactor dictating the contribution of each loss on the gradient, we make use of the work by Kendall et al. [69] that weighs the different losses based on their homoscedastic uncertainty, thereby reducing the number of hyperparameters. The model is trained using the SGD optimizer, a momentum of 0.95 and a learning rate of  $2 \times 10^{-3}$ . The batch size is set to 5, and the training runs for 70 epochs for UAVDT and 150 epochs for VisDrone2019. In the cases where not all the cells inside the bounding box are used to localize the center, a region defined as an "ignore" region, with a size set to 20% that of the object in question, is created around the used cells.

### 4.3 Testing Details

During testing, the image is rescaled to have the width and the height as multiples of 32, a requirement based on the architecture of HRNet. The confidence threshold for selecting centers is 0.05, similar to previous methods. We set the  $\rho$  to 10 for the UAVDT Dataset and 5 for Visdrone. The latter requires a smaller  $\rho$  because it contains small objects. Once the butterfly fields are extracted, we use the voting mechanism introduced in Section 3.3 to aggregate and obtain the predicted bounding boxes. The resulting bounding boxes are then refined using soft non-maximum suppression [77].

Both datasets are evaluated using average precision (AP). UAVDT benchmark requires the AP to be calculated at a specific IOU threshold of 0.7. On the contrary, VisDrone2019 requires an evaluation protocol similar to MS-COCO [43] where the average precision as well as average recall are calculated at different IOU thresholds. These metrics are shown in Table 2 and are explained in detail in [43].

## 5 Results

### 5.1 Comparison to State-of-the-art

Our method is able to significantly outperform previous state-of-the-art methods on UAVDT dataset. As seen in Table 1 and comparing it to published state-of-the-art method (NDFT F-RCNN), our butterfly detector reaches an average precision of 58.1%, exceeding it by  $\sim 10\%$ . We also validate our motivation for the use of fields in aerial images. This is shown by the drastic improvement of  $\sim 14\%$  in average precision on medium to high altitude images as well as an improvement of  $\sim 17\%$  on images taken from bird view. We also train a state-of-the-art anchor-free object detector, CenterNet, on UAVDT. It uses an Hourglass architecture with input size  $512 \times 512$ . Testing is done on the full resolution. Our method outperforms CenterNet by  $\sim 4\%$  (Table 1) while obtaining

Table 1: Comparison of AP (%), False Positives (FP), and Recall (%) with state-of-the-art methods on UAVDT datasets. We outperform previous methods on medium to high altitude images as well as bird-view images. These are significant challenges that general object detectors face on aerial images.  $BD_{NV}$  is our method with no voting mechanism and using all predicted bounding boxes. The subscript refers to the number of feature cells used to localize the center.  $BD_{full}$  makes use of all feature cell inside the ground truth bounding box. \*: trained and tested by us. *fps*: frames per second

	Backbone	FP	Recall	<i>avg.</i>	Flying Altitude			Camera View			Time		
					<i>low</i>	<i>medium</i>	<i>high</i>	<i>front</i>	<i>side</i>	<i>bird</i>	<i>day</i>	<i>night</i>	<i>fps</i>
YOLOv3[15]	DarkNet-53 +FF	–	–	20.6	–	–	–	–	–	–	–	–	–
SSD [23]	VGG-16	2,111k	49.28	33.6	–	–	–	–	–	–	–	–	42
CADNet [55]	VGG-16 +FF	–	–	43.6	–	–	–	–	–	–	–	–	–
Faster R-CNN [17]	ResNet-101	–	–	45.64	68.14	49.71	18.70	53.34	68.02	27.05	45.63	52.14	2.8
GANet [54]	ResNet-50 +FF	–	–	47.2	–	–	–	–	–	–	–	–	–
NDFT F-RCNN [16]	ResNet-101	–	–	47.91	<b>74.84</b>	56.24	20.55	<b>64.88</b>	<b>67.50</b>	28.79	45.91	64.16	–
CenterNet* [47]	Hourglass-104	1,204k	66.0	53.95	70.21	60.22	27.65	57.13	66.29	38.85	62.63	67.24	6
$BD_{16}$	ResNet-50	453k	62.6	52.74	68.32	61.87	21.44	57.44	65.15	34.55	62.58	69.27	12.2
$BD_{NV}$	HRNet-W32	2,881k	<b>71.28</b>	45.02	58.21	51.79	21.52	48.83	50.62	32	50.29	64.61	–
$BD_1$	HRNet-W32	412k	67.46	51.44	63.19	56.3	29.57	49.91	62.13	42.14	58.07	62.94	18.4
$BD_{full}$	HRNet-W32	224k	63.74	54.28	67.6	59.77	31.26	55.16	63.88	42.79	60.84	70.88	15.4
$BD_{16}$	HRNet-W32	313k	69.46	<b>58.1</b>	67.84	<b>61.93</b>	<b>40.6</b>	56.28	66.81	<b>51.57</b>	<b>62.97</b>	<b>69.59</b>	<b>18.4</b>

Table 2: Comparison of AP (average precision), AR (average recall), True Positives (TP), and False Positives (FP) with state-of-the-art methods on VisDrone dataset. As shown, we outperform previous methods on all metrics. *det.*: specifies the maximum detections chosen. \*\*: trained and tested by us. \*: use of sub-pixel localization

	AP(% , det.=500)			AR(% , IoU <sub>0.5:0.95</sub> )				TP	FP
	IoU <sub>0.5:0.95</sub>	IoU <sub>0.5</sub>	IoU <sub>0.75</sub>	<i>det.</i> = 1	<i>det.</i> = 10	<i>det.</i> = 100	<i>det.</i> = 500	IoU <sub>0.75</sub>	IoU <sub>0.75</sub>
FRCNN[17] +FPN[78]	21.4	40.7	19.9	–	–	–	–	–	–
CenterNet [11]*	22.96	42.05	21.81	0.68	6.14	32.56	32.56	15.9k	36k
ClusDet [57]	28.4	53.2	26.4	–	–	–	–	–	–
BD	28.93	54.92	26.15	0.98	7.06	36.1	40.93	17.8k	46.7k
<b>BD*</b>	<b>30.15</b>	<b>54.67</b>	<b>28.66</b>	<b>0.98</b>	<b>7.11</b>	<b>37.44</b>	<b>41.41</b>	19.3k	31.7k

similar results using ResNet-50. Compared to CenterNet, we especially achieve better results on medium and high altitude images (+1.7% and +13% respectively) as well as on bird-view images ( $\sim +13\%$ )

We also obtain state-of-the-art results on the VisDrone2019 dataset (Table 2). We report our results on the validation set. Our detector is able to attain better performance than all previous methods. ClusDet [57], a two-stage detector, achieves good results since it runs inference on multiple small patches of the image as well as the full image. Although this allows it to be more fine-grained, we outperform it by  $\sim 2.2\%$ . Comparing it to CenterNet [12], an anchor-free object detector that suffer from the limitations discussed previously, our method outperforms it by  $\sim 7\%$  when averaging over

the different IoUs and by  $\sim 12\%$  at IoU=0.5. The improvement of the AP at IoU=0.75 indicates that our method is able to correctly predict the scale of the objects.

## 5.2 Inference Speed

We also focus on the inference time of previous methods and compare them to our butterfly detector (Table 1). Most methods did not provide their inference time for comparison. As observed, Faster R-CNN has a low fps due to it being a two-stage detector that relies on a region-proposal network(RPN). Since the architecture and base network of NDFT F-RCNN is based on Faster R-CNN, it is expected that it also does not run in real-time reaching around  $3\text{ fps}$ . This is also the case for CADNet that also makes use of a more complicated RPN. Our method performs better than all previous state-of-the-art while running at a speed of 18.4 fps, a real-time performance. All reported inference speeds were obtained using a GTX 1080 Ti.

## 6 Ablation Study

**Effect of Fields.** In order to showcase the benefits of fields, we compare in Table 1 our method ( $BD_{16}$ ) to the same method but using only one feature cell responsible for the center ( $BD_1$ ). As can be observed, a gain of  $\sim 7\%$  is achieved when different locations in a field are used to vote for the localization of the center instead of only using the center cell. We observe that we achieve a boost of  $2\%$  in recall as well. This is due to the fact that different bounding boxes are being proposed based on the extracted features at different output cells.

**Effect of Voting Mechanism.** To study the advantages of voting, we measure the performance of the model when all the bounding boxes predicted by the activated output feature cells are used. Our method without the voting mechanism achieves an AP of 45.02% (Table 1) which is almost on-par with Faster R-CNN, a slower two stage detector. As expected, we observe a large difference between true positives and false positives leading to a high recall but low average precision. When the voting mechanism described in section 3.3 is applied, a significant decrease in the false positives is achieved while preserving the recall. As a result, a significant boost in the AP is observed of around 13% increasing from 45.02% to 58.1%. This improvement is also shown over all the different attributes of the UAVDT dataset.

**Effect of Backbone.** In order to better verify the positive effect of our method, we replace the HRNet-W32 with a commonly-used backbone, ResNet-50. Our method with a ResNet-50 architecture is able to outperform previous methods while achieving similar results to CenterNet. These previous methods make use of larger architectures as well as fusion between different feature scales (Table 1). HRNet-W32 is used in our method since it allows learning from different feature scales similar to CenterNet’s backbone, Hourglass-104.

**Number of Localizing Cells in a Field.** We also study the effect of the number of points in a field localizing the center of an object on the performance. We evaluate the performance of our method while using one feature cell which is the center ( $BD_1$ ), 16 feature cells around the center ( $BD_{16}$ ), and all feature cells lying inside the bounding box ( $BD_{full}$ ). Table 1 shows the numbers of true and false positives as well as the average precision per attribute for these three configurations. We notice



(a) UAVDT

(b) VisDrone2019

Figure 5: Qualitative results of the Butterfly Detector on UAVDT and VisDrone dataset. Our method was able to perfectly detect partly occluded objects (*e.g.*, images from UAVDT). This is because our Butterfly Detector leverages several locations in a field pointing to the center through a voting mechanism

that the best settings is using 16 cells in a 4x4 window to localize the center and predict the width and the height of the bounding box.  $BD_{full}$  works by utilizing all features inside the bounding box. Since ground truth boxes are axis-aligned, a bounding box of a car moving diagonally relative to the image will contain areas of the road. In such cases, using the full boxes requires classifying some roads with high confidence leading to a difficult training process. Thus, this can explain the slightly lower performance compared to  $BD_{16}$ . Nonetheless, our detector with the full bounding box still outperforms previous methods (Table 1). As for  $BD_1$ , it makes use of only the center feature to localize the object which makes it challenging to train and decode as explained previously.

**Effect of Sub-pixel Localization.** To better observe the effect of the sub-pixel localization, we perform this ablation study on VisDrone2019. As observed in Table 2, using the sub-pixel technique improves the performance of the model by  $\sim 1.2\%$ . This method also helped in drastically reducing the false positives while also increasing the true positives.

## 7 Conclusions

In this paper, we propose a new object detector, Butterfly Detector, that utilizes composite fields over the spatial feature map to localize the center of objects and predict their sizes. Using the extracted fields for an object center, we apply a voting mechanism to further improve our detector especially in cases of occlusion or viewing angle variations. We evaluated our approach on UAVDT and

VisDrone and show that it outperforms the previous state-of-the-art methods. We further perform an ablation study to demonstrate the benefits of butterfly fields, coupled with a voting process and sub-pixel localization, for object detection in aerial images.

## References

- [1] Meticulous Market Research Pvt. Ltd. Smart cities market worth \$545.7 billion by 2025, growing at a cagr of 22.9% from 2019- global market opportunity analysis and industry forecasts by meticulous research®. "GlobeNewswire", Jun 2020.
- [2] Emmanouil Barmponakis and Nikolas Geroliminis. On the new era of urban traffic monitoring with massive drone data: The pneuma large-scale field experiment. *Transportation Research Part C: Emerging Technologies*, 111:50–71, 2020.
- [3] Kevin Crowston. Amazon mechanical turk: A research tool for organizations and information systems scholars. In Anol Bhattacharjee and Brian Fitzgerald, editors, *Shaping the Future of ICT Research. Methods and Approaches*, pages 210–221, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg.
- [4] Zheng Wang and Jih-Biing Sheu. Vehicle routing problem with drones. *Transportation research part B: methodological*, 122:350–364, 2019.
- [5] Chun Cheng, Yossiri Adulyasak, and Louis-Martin Rousseau. Drone routing with energy function: Formulation and exact algorithm. *Transportation Research Part B: Methodological*, 139:364–387, 2020.
- [6] Pedro L Gonzalez-R, David Canca, Jose L Andrade-Pineda, Marcos Calle, and Jose M Leon-Blanco. Truck-drone team logistics: A heuristic approach to multi-drop route planning. *Transportation Research Part C: Emerging Technologies*, 114:657–680, 2020.
- [7] Mohamed Salama and Sharan Srinivas. Joint optimization of customer location clustering and drone-based routing for last-mile deliveries. *Transportation Research Part C: Emerging Technologies*, 114:620 – 642, 2020.
- [8] Chase C. Murray and Ritwik Raj. The multiple flying sidekicks traveling salesman problem: Parcel delivery with multiple drones. *Transportation Research Part C: Emerging Technologies*, 110:368 – 398, 2020.
- [9] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2009.
- [10] Qi Fan, Wei Zhuo, Chi-Keung Tang, and Yu-Wing Tai. Few-shot object detection with attention-rpn and multi-relation detector. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [11] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018.

- [12] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv*, 2019.
- [13] 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.
- [14] Pengfei Zhu, Longyin Wen, Xiao Bian, Ling Haibin, and Qinghua Hu. Vision meets drones: A challenge. *arXiv*, 2018.
- [15] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016.
- [16] 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. *arXiv*, 2019.
- [17] 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.
- [18] Jian Ding, Nan Xue, Yang Long, Gui-Song Xia, and Qikai Lu. Learning roi transformer for oriented object detection in aerial images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [19] Zhe Cao, Gines Hidalgo, Tomas Simon, Shih-En Wei, and Yaser Sheikh. Openpose: realtime multi-person 2d pose estimation using part affinity fields. *arXiv*, 2018.
- [20] Sven Kreiss, Lorenzo Bertoni, and Alexandre Alahi. Pifpaf: Composite fields for human pose estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11977–11986, 2019.
- [21] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–286, 2018.
- [22] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [23] 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 (ECCV)*, pages 21–37. Springer, 2016.
- [24] Gregory P. Meyer, Ankit Laddha, Eric Kee, Carlos Vallespi-Gonzalez, and Carl K. Wellington. Lasernet: An efficient probabilistic 3d object detector for autonomous driving. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

- [25] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [26] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [27] Lijie Liu, Jiwen Lu, Chunjing Xu, Qi Tian, and Jie Zhou. Deep fitting degree scoring network for monocular 3d object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [28] Surendra Gupte, Osama Masoud, Robert FK Martin, and Nikolaos P Papanikolopoulos. Detection and classification of vehicles. *IEEE Transactions on intelligent transportation systems*, 3(1):37–47, 2002.
- [29] D Chaudhuri, NK Kushwaha, and Ashok Samal. Semi-automated road detection from high resolution satellite images by directional morphological enhancement and segmentation techniques. *IEEE journal of selected topics in applied earth observations and remote sensing*, 5(5):1538–1544, 2012.
- [30] Jonathan Weber and Sébastien Lefèvre. Spatial and spectral morphological template matching. *Image and Vision Computing*, 30(12):934–945, 2012.
- [31] Yudong Lin, Hongjie He, Zhongke Yin, and Fan Chen. Rotation-invariant object detection in remote sensing images based on radial-gradient angle. *IEEE Geoscience and Remote Sensing Letters*, 12(4):746–750, 2014.
- [32] Salman Ahmadi, MJ Valadan Zoej, Hamid Ebadi, Hamid Abrishami Moghaddam, and Ali Mohammadzadeh. Automatic urban building boundary extraction from high resolution aerial images using an innovative model of active contours. *International Journal of Applied Earth Observation and Geoinformation*, 12(3):150–157, 2010.
- [33] 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.
- [34] Kun Liu, Henrik Skibbe, Thorsten Schmidt, Thomas Blein, Klaus Palme, Thomas Brox, and Olaf Ronneberger. Rotation-invariant hog descriptors using fourier analysis in polar and spherical coordinates. *International Journal of Computer Vision*, 106(3):342–364, 2014.
- [35] Jon Almazán, Alicia Fornés, and Ernest Valveny. Deformable hog-based shape descriptor. In *2013 12th International Conference on Document Analysis and Recognition*, pages 1022–1026. IEEE, 2013.
- [36] Gong Cheng, Junwei Han, Lei Guo, Xiaoliang Qian, Peicheng Zhou, Xiwen Yao, and Xintao Hu. Object detection in remote sensing imagery using a discriminatively trained mixture model. *ISPRS Journal of Photogrammetry and Remote Sensing*, 85:32–43, 2013.

- [37] Gong Cheng, Junwei Han, Peicheng Zhou, and Lei Guo. Multi-class geospatial object detection and geographic image classification based on collection of part detectors. *ISPRS Journal of Photogrammetry and Remote Sensing*, 98:119–132, 2014.
- [38] Line Eikvil, Lars Aurdal, and Hans Koren. Classification-based vehicle detection in high-resolution satellite images. *ISPRS Journal of Photogrammetry and Remote Sensing*, 64(1):65–72, 2009.
- [39] Helmut Grabner, Thuy Thi Nguyen, Barbara Gruber, and Horst Bischof. On-line boosting-based car detection from aerial images. *ISPRS Journal of Photogrammetry and Remote Sensing*, 63(3):382–396, 2008.
- [40] Wanceng Zhang, Xian Sun, Hongqi Wang, and Kun Fu. A generic discriminative part-based model for geospatial object detection in optical remote sensing images. *ISPRS journal of photogrammetry and remote sensing*, 99:30–44, 2015.
- [41] Yuxiang Zhang, Liangpei Zhang, Bo Du, and Shugen Wang. A nonlinear sparse representation-based binary hypothesis model for hyperspectral target detection. *IEEE Journal of selected topics in applied earth observations and remote sensing*, 8(6):2513–2522, 2014.
- [42] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [43] 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.
- [44] 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.
- [45] Jifeng Dai, Yi Li, Kaiming He, and Jian Sun. R-fcn: Object detection via region-based fully convolutional networks. In *Advances in neural information processing systems*, pages 379–387, 2016.
- [46] Lars Sommer, Lucas Steinmann, Arne Schumann, and Jürgen Beyerer. Systematic evaluation of deep learning based detection frameworks for aerial imagery. In Firooz A. Sadjadi and Abhijit Mahalanobis, editors, *Automatic Target Recognition XXVIII*, volume 10648, pages 1 – 13. International Society for Optics and Photonics, SPIE, 2018.
- [47] Lars Wilko Sommer, Tobias Schuchert, and Jürgen Beyerer. Fast deep vehicle detection in aerial images. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 311–319. IEEE, 2017.
- [48] Jinze Li, Rong Wang, and Jianwei Ding. Tiny vehicle detection from uav imagery. In Yongtian Wang, Qingmin Huang, and Yuxin Peng, editors, *Image and Graphics Technologies and Applications*, pages 370–381, Singapore, 2019. Springer Singapore.

- [49] Lars Sommer, Arne Schumann, Tobias Schuchert, and Jurgen Beyerer. Multi feature deconvolutional faster r-cnn for precise vehicle detection in aerial imagery. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 635–642. IEEE, 2018.
- [50] Michael Ying Yang, Wentong Liao, Xinbo Li, and Bodo Rosenhahn. Deep learning for vehicle detection in aerial images. In *2018 25th IEEE International Conference on Image Processing (ICIP)*, pages 3079–3083. IEEE, 2018.
- [51] Tobias Ringwald, Lars Sommer, Arne Schumann, Jurgen Beyerer, and Rainer Stiefelhagen. Uav-net: A fast aerial vehicle detector for mobile platforms. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- [52] Tianyu Tang, Shilin Zhou, Zhipeng Deng, Huanxin Zou, and Lin Lei. Vehicle detection in aerial images based on region convolutional neural networks and hard negative example mining. *Sensors*, 17(2):336, 2017.
- [53] Zhipeng Deng, Hao Sun, Shilin Zhou, Juanping Zhao, and Huanxin Zou. Toward fast and accurate vehicle detection in aerial images using coupled region-based convolutional neural networks. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 10(8):3652–3664, 2017.
- [54] Yuanqiang Cai, Dawei Du, Libo Zhang, Longyin Wen, Weiqiang Wang, Yanjun Wu, and Siwei Lyu. Guided attention network for object detection and counting on drones. *arXiv*, 2019.
- [55] Kaiwen Duan, Dawei Du, Honggang Qi, and Qingming Huang. Detecting small objects using a channel-aware deconvolutional network. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [56] Wei Li, Hongliang Li, Qingbo Wu, Xiaoyu Chen, and King Ngi Ngan. Simultaneously detecting and counting dense vehicles from drone images. *IEEE Transactions on Industrial Electronics*, 2019.
- [57] Fan Yang, Heng Fan, Peng Chu, Erik Blasch, and Haibin Ling. Clustered object detection in aerial images. *CoRR*, abs/1904.08008, 2019.
- [58] Hao Long, Yinung Chung, Zhenbao Liu, and Shuhui Bu. Object detection in aerial images using feature fusion deep networks. *IEEE Access*, 7:30980–30990, 2019.
- [59] Bodi Ma, Zhenbao Liu, Feihong Jiang, Yuehao Yan, Jinbiao Yuan, and Shuhui Bu. Vehicle detection in aerial images using rotation-invariant cascaded forest. *IEEE Access*, 7:59613–59623, 2019.
- [60] Lachlan Tychsen-Smith and Lars Petersson. Denet: Scalable real-time object detection with directed sparse sampling. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 428–436, 2017.
- [61] Xinggong Wang, Kaibing Chen, Zilong Huang, Cong Yao, and Wenyu Liu. Point linking network for object detection. *arXiv*, 2017.

- [62] Hei Law, Yun Teng, Olga Russakovsky, and Jia Deng. Cornernet-lite: Efficient keypoint based object detection. *arXiv*, 2019.
- [63] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019.
- [64] Lichao Huang, Yi Yang, Yafeng Deng, and Yinan Yu. Densebox: Unifying landmark localization with end to end object detection. *CoRR*, abs/1509.04874, 2015.
- [65] Jiaqi Wang, Kai Chen, Shuo Yang, Chen Change Loy, and Dahua Lin. Region proposal by guided anchoring. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [66] Chenchen Zhu, Yihui He, and Marios Savvides. Feature selective anchor-free module for single-shot object detection. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [67] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [68] Tao Kong, Fuchun Sun, Huaping Liu, Yuning Jiang, and Jianbo Shi. Foveabox: Beyond anchor-based object detector. *arXiv*, 2019.
- [69] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [70] J. Uhrig, E. Rehder, B. Fröhlich, U. Franke, and T. Brox. Box2pix: Single-shot instance segmentation by assigning pixels to object boxes. In *2018 IEEE Intelligent Vehicles Symposium (IV)*, pages 292–299, 2018.
- [71] Davy Neven, Bert De Brabandere, Marc Proesmans, and Luc Van Gool. Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [72] Bowen Cheng, Maxwell D. Collins, Yukun Zhu, Ting Liu, Thomas S. Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [73] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2019.
- [74] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in neural information processing systems*, pages 5574–5584, 2017.

- [75] Lorenzo Bertoni, Sven Kreiss, and Alexandre Alahi. Monoloco: Monocular 3d pedestrian localization and uncertainty estimation. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019.
- [76] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *CoRR*, abs/1708.02002, 2017.
- [77] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S. Davis. Soft-nms – improving object detection with one line of code. In *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [78] 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.