

Drone-based Object Counting by Spatially Regularized Regional Proposal Network

Meng-Ru Hsieh¹, Yen-Liang Lin², and Winston H. Hsu¹

¹National Taiwan University, Taipei, Taiwan ²GE Global Research, Niskayuna, NY, USA mrulafi@gmail.com, yenlianglintw@gmail.com, whsu@ntu.edu.tw

Abstract

Existing counting methods often adopt regression-based approaches and cannot precisely localize the target objects, which hinders the further analysis (e.g., high-level understanding and fine-grained classification). In addition, most of prior work mainly focus on counting objects in static environments with fixed cameras. Motivated by the advent of unmanned flying vehicles (i.e., drones), we are interested in detecting and counting objects in such dynamic environments. We propose Layout Proposal Networks (LPNs) and spatial kernels to simultaneously count and localize target objects (e.g., cars) in videos recorded by the drone. Different from the conventional region proposal methods, we leverage the spatial layout information (e.g., cars often park regularly) and introduce these spatially regularized constraints into our network to improve the localization accuracy. To evaluate our counting method, we present a new large-scale car parking lot dataset (CARPK) that contains nearly 90,000 cars captured from different parking lots. To the best of our knowledge, it is the first and the largest drone view dataset that supports object counting, and provides the bounding box annotations.

1. Introduction

With the advent of unmanned flying vehicles, new potential applications emerge for unconstrained images and videos analysis for aerial view cameras. In this work, we address the counting problem for evaluating the number of objects (e.g., cars) in drone-based videos. Prior methods [10, 2, 1] for monitoring the parking lot often assume that the locations of the monitored objects of a scene are already known in advance and the cameras are fixed, and cast car counting as a classification problem, which makes conventional car counting methods not directly applicable in unconstrained drone videos.



Figure 1. We propose a Layout Proposal Network (LPN) to localize and count objects in drone videos. We introduce the spatial constraints for learning our network to improve the localization accuracy. Detailed network structure is shown in Figure 4.

Current object counting methods often learn a regression model that maps the high-dimensional image space into non-negative counting numbers [30, 18]. However, these methods can not generate precise object positions, which limits the further investigation and applications (e.g., recognition).

We observe that there exists certain layout patterns for a group of object instances, which can be utilized to improve the object counting accuracy. For example, cars are often parked in a row and animals are gathered in a certain layout (e.g., fish torus and duck swirl). In this paper, we introduce a novel Layout Proposal Network (LPN) that counts and localizes objects in drone videos (Figure 1). Different from existing object proposal methods, we introduce a new spatially regularized loss for learning our Layout Proposal Network. Note that our method learns the general adjacent relationship between object proposals and is not specific to a certain scene.

Our spatially regularized loss is a weighting scheme that

Table 1. Comparison of aerial view car-related datasets. In contrast to the PUCPR dataset, our dataset supports a counting task with bounding box annotations for all cars in a single scene. Most important of all, compared to other car datasets, our CARPK is the only dataset in drone-based scenes and also has a large enough number in order to provide sufficient training samples for deep learning models.

Dataset	Sensor	Multi Scenes	Resolution	Annotation Format	Car Numbers	Counting Support
OIRDS [28]	satellite	\checkmark	low	bounding box	180	\checkmark
VEDAI [20]	satellite	\checkmark	low	bounding box	2,950	\checkmark
COWC [18]	aerial	\checkmark	low	car center point	32,716	\checkmark
PUCPR [10]	camera	×	high	bounding box	192,216	×
CARPK [ours]	drone	\checkmark	high	bounding box	89,777	\checkmark



(d)

(e)

Figure 2. (a), (b), (c), (d), and (e) are the example scenes of OIRDS [28], VEDAI [20], COWC [18], PUCPR [10], and CARPK (ours) dataset respectively (two images for each dataset). Comparing to (a), (b), and (c), the PUCPR dataset and the CARPK dataset have greater number of cars in a single scene which is more appropriate for evaluating the counting task.

re-weights the importance scores for different object proposals and encourages region proposals to be placed in correct locations. It can also generally be embedded in any object detector system for object counting and detection. By exploiting spatial layout information, we improve the average recall of state-of-the-art region proposal approaches on a public PUCPR dataset [10] (from 59.9% to 62.5%).

For evaluating the effectiveness and reliability of our approach, we introduce a new large-scale counting dataset CARPK (Table 1). Our dataset contains 89,777 cars, and provides bounding box annotations for each car. Also, we consider the sub-dataset PUCPR of PKLot [10] which is the one that the scenes are closed to the aerial view in the PKLot dataset. Instead of a fixed camera view from a high story building (Figure 2) in the PUCPR dataset, our new CARPK dataset provide the first and the largest-scale drone view parking lot dataset in unconstrained scenes. Besides, the PUCPR dataset can be only used in conjunction with a classification task, which classifies the pre-cropped images (car or not car) with given locations. Moreover, the PUCPR dataset only annotates partial region-of-interest parking ar-

eas, and is therefore unable to support a counting task. Since our task is to count objects in images, we also annotate all cars in single full-image for the partial PUCPR dataset. The contents of our CARPK dataset are unscripted and diverse in various scenes for 4 different parking lots. To the best of our knowledge, our dataset is the first and the largest drone-based dataset that can support a counting task with manually labelled annotations for numerous cars in full images. The main contributions of this paper are summarized as follows:

- 1. To our knowledge, this is the first work that leverages spatial layout information for object region proposal. We improve the average recall of the state-of-the-art region proposal methods (i.e., 59.9% [22] to 62.5%) on a public PUCPR dataset.
- 2. We introduce a new large-scale car parking lot dataset (CARPK) that contains nearly 90,000 cars in dronebased high resolution images recorded from the diverse scenes of parking lots. Most important of all, compared to other parking lot datasets, our CARPK

dataset is the first and the largest dataset of parking lots that can support counting¹.

3. We provide in-depth analyses for different decision choices of our region proposal method, and demonstrate that utilizing layout information can considerably reduce the proposals and improve the counting results.

2. Related Work

2.1. Object Counting

Most contemporary counting methods can be broadly divided into two categories. One is counting by regression method, the other is counting by detection instance [17, 14]. Regression counters are usually a mapping of the high-dimension image space into non-negative counting numbers. Several methods [3, 6, 7, 8, 15] try to predict counts by using global regressors trained with low-level features. However, global regression methods ignore some constraints, such as the fact that people usually walk on the pavement and the size of instances. There are also a number of density regression-based methods [23, 16, 5] which can estimate object counts by the density of a countable object and then aggregate over that density.

Recently, a wealth of works introduce deep learning into the crowd counting task. Instead of counting objects for constrained scenes in the preivous works, Zhang et al. [31] address the problem of cross-scene crowd counting task, which is the weakness of the density estimation method in the past. Sindagi et al. [26] incorporate global and local contextual information for better estimating the crowd counts. Mundhenk et al. [18] evaluate the number of cars in a subspace of aerial imagery by extracting representations of image patches to approximate the appearance of object groups. Zhang et al. [32] leverage FCN and LSTM to jointly estimate the vehicle density and counts in low resolution videos from city cameras. However, the regressionbased methods can not generate precise object positions, which seriously limits the further investigation and application (e.g., high-level understanding and fine-grained classification).

2.2. Object Proposals

Recent years have seen deep networks for region proposals developing well. Because detecting objects at several positions and scales during inference time requires a computationally demanding classifier, the best way to solve this problem is to look at a tiny subset of possible positions. A number of recent works prove that deep networks-based region proposal methods have surpassed the previous works [29, 4, 33, 9], which are based on the low-level cues, by a large margin.

DeepMask [19], which is developed for learning segmentation proposals, has, compared to Selective Search [29], ten times fewer proposals (100 v.s. 1000) at the same performance. The state-of-the-art object proposal method, Region Proposal Networks (RPNs) [22], has also shown that they just need 300 proposals and can surpass the result of 2000 proposals generated by [29]. Other works like Multibox [27] and Deepbox [33] also have higher proposal recall with fewer number of region proposals than the previous works which are based on low-level cues. However, none of these region proposal methods have considered the spatial layout or the relation between recurring objects. Hence, we propose a Layout Proposal Networks (LPNs) that leverages thus structure information to achieve higher recall while using a smaller number of proposals.

3. Dataset

Since there is a lack of large standardized public datasets that contain numerous collections of cars in drone-based images, it has been difficult to create an automated counting system for deep learning models. For instance, OIRDS [28] has merely 180 unique cars. The recent car-related dataset VEDAI [20] has 2,950 cars, but these are still too few to utilize for the deep learners. A newer dataset COWC [18] has 32,716 cars, but the resolutions of images remain low. It has only 24 to 48 pixels per car. Besides, rather than labelling in the format of bounding box, the annotation format is the center pixel point of a car which can not support further investigation, such as car model retrieval, statistics of brands of car, and exploring which kind of car most people will drive in the local area. Moreover, all above datasets are low resolution images and cannot provide detail informations for learning a fine-grained deep model. The problems of existing dataset are : 1) low resolution images which might harm the performance of the model trained on them and 2) less car numbers in the dataset which has the potential to cause overfitting during training a deep model.

Because existing datasets have these aforementioned problems, we have created a large-scale car parking lot dataset from drone view images, which are more appropriate to deep learning algorithms. It supports object counting, object localizing, and further investigations by providing the annotations in terms of bounding boxes. The most similar public dataset to ours, which also has the high resolution of car images, is the sub-dataset PUCPR of PKLot [10], which provides a view from the 10th floor of a building and therefore similar to drone view images to a certain degree. However, the PUCPR dataset can be only used in conjunction with a classification task, which classifies the precropped images (car or not car) with given locations. Moreover, this dataset has only annotated a portion of cars (100

¹The images and annotations of our CARPK and PUCPR+ are available at https://lafi.github.io/LPN/



Figure 3. The key idea of our spatial layout scores. A predicted position which has more nearby cars can get higher confidence scores and has higher probability to be the position where the car is.

certain parking spaces) from total 331 parking spaces in a single image, making it unable to support both counting and localizing tasks. Hence, we complete the annotations for all cars in a single image from the partial PUCPR dataset, called PUCPR+ dataset, which now has nearly 17,000 cars in total. Besides the incomplete annotation problem of the PUCPR, it has a fatal issue that their camera sensors are fixed and set in the same place, making the image scene of dataset completely the same – causing the deep learning model to encounter a dataset bias problem.

For this reason, we introduce a brand new dataset CARPK that the contents of our dataset are unscripted and diverse in various scenes for 4 different parking lots. Our dataset also contains approximately 90,000 cars in total with the view of drone. It is different from the view of camera from high story building in the PUCPR dataset. This is a large-scale dataset for car counting in the scenes of diverse parking lots. The image set is annotated by providing a bounding box per car. All labeled bounding boxes have been well recorded with the top-left points and the bottomright points. Cars located on the edge of the image are included as long as the marked region can be recognized and it is sure that the instance is a car. To the best of our knowledge, our dataset is the first and the largest drone viewbased parking lot dataset that can support counting with manually labeled annotations for a great amount of cars in a full-image. The details of dataset are listed in Table 1 and some examples are shown in Figure 2.

4. Method

Our object counting system employs a region proposal module which takes regularized layout structure into account. It is a deep fully convolutional network that takes an image of arbitrary size as the input, and outputs the objectagnostic proposals which likely contain the instance. The entire system is a single unified framework for object counting (Figure 1). By leveraging the spatial information of the object of recurring instances, LPNs module is not only concerning the possible positions but also suggesting the object detection module which direction it should look at in the image.

4.1. Layout Proposal Network

We observed that there exists certain layout patterns for a group of object instances, which can be used to predict objects that might appear adjacently in the same direction or near the same instances. Hence, we design a novel region proposal module that can leverage the structure layout and gather the confidence scores from nearby objects in certain directions (Figure 3).

We comprehensively describe the designed network structure of LPNs (Figure 4) as follows. Similar to RPNs [22], the network generates region proposals by sliding a small network over the shared convolutional feature map. It takes as input an 3×3 windows on last convolutional layer for reducing the representation dimensions, and then feeds features into two sibling 1×1 convolutional layers, where one is for localization and the other is for classifying whether the box belongs to foreground or background. The difference is that our loss function introduces the spatially regularized weights for the predicted boxes at each location. With the weights from spatial information, we minimize the loss of multi-task object function in networks. The loss function we use on each image is defined as:

$$L(\{u_i\}, \{q_i\}, \{p_i\}) = \frac{1}{N_{fg}} \sum_i K(c_i, N_i^*; u_i^*) \cdot L_{fg}(u_i, u_i^*) + \gamma \frac{1}{N_{bg}} \sum_i L_{bg}(q_i, q_i^*) + \lambda \frac{1}{N_{loc}} \sum_i L_{loc}(u_i^*, p_i, g_i^*)$$
(1)

where N_{fg} and N_{bg} are the normalized terms of the number matching default boxes for foreground and background. N_{loc} is the same as N_{fg} in that it only considers the number of foreground classes. The default box is marked $u_i^* = 1$ if the default box has an Intersection-over-Union (IoU) overlap higher than 0.7 with the ground-truth box, or the default box which has the highest IoU overlap with a groundtruth box; otherwise, it is marked $q_i^* = 0$ if the IoU overlap is lower than 0.3. The $L_{fg}(u_i, u_i^*) = -log[u_i u_i^*]$ and $L_{bg}(q_i, q_i^*) = -log[(1 - q_i)(1 - q_i^*)]$ are the negative loglikelihood that we want to minimize for true classes. Here, the *i* is the index of predicted boxes. In front of the foreground loss, *K* represents that we apply the spatially regularized weights for re-weighting the objective score of each predicted box. The weight is obtained by a Gaussian spatial



Figure 4. The structure of the Layout Proposal Networks. At the loss layer, the structure weights are integrated for re-weighting the candidates to have better structure proposals. See more details in Section 4.2.

kernel for the center position c_i of predicted box. It will give a rearranged weight according to the m neighbor groundtruth box centers, which are near to the c_i . The real neighbor centers for c_i are denoted as $N_i^* = \{c_1^*, ..., c_m^*\} \in S^{c_i}$, which fall inside the spatial window pixels size S on the input image. We use S = 255 in this paper to obtain a larger spatial range.

The L_{loc} is the localization loss, which is a robust loss function [11]. This term is only active for foreground predicted boxes $(u_i^* = 1)$, otherwise 0. Similar to [22], we calculate the loss of offsets between the foreground predicted box p_i and the ground truth box g_i with their center position (x, y), width (w), and height (h) based on the default box (d).

$$L_{loc}(u_{i}^{*}, p_{i}, g_{i}^{*}) = \sum_{i \in fg} \sum_{v \in \{x, y, w, h\}} u_{i}^{*} smooth_{L1}(p_{i}^{v}, g_{i}^{v*})$$
(2)

, where q_i^{v*} (similar to p_i^v) is defined as below:

$$g_i^{x*} = (g_i^x - d_i^x)/d_i^w, \qquad g_i^{y*} = (g_i^y - d_i^y)/d_i^h$$

$$g_i^{w*} = \log(g_i^w/d_i^w), \qquad g_i^{h*} = \log(g_i^h/d_i^h)$$
(3)

In our experiment, we set γ and λ to be 1. Besides, in order to handle the small objects, instead of conv5-3 layer, we select conv4-3 layer features for obtaining better tiling default box stride on the input image and choose default box sizes approximately four times smaller ($16 \times 16, 40 \times 40, 100 \times$ 100) than the default setting ($128 \times 128, 256 \times 256, 512 \times$ 512).

4.2. Spatial Pattern Score

Most of the objects of an instance exhibit a certain pattern between each other. For instance, cars will align in one direction on a parking lot and ships will hug the shore regularly. Even in biology, we can also find collective animal behavior that makes them look into a certain layout (e.g., fish torus, duck swirl, and ant mill). Hence, we introduce a method for re-weighting the region proposals in the training phase in an end-to-end manner. The proposed method can reduce the number of proposals in the inference phase for abating the computational cost of the counting and detection processes. It is especially important on embedded devices, such as the drone, to lower power consumption as the battery power only can provide the drone with energy to fly a mere 20 minutes.

For designing the pattern of layout, we apply different direction 2D Gaussian spatial kernels K (see Eq. 1) on the space of input images, where the center of the Gaussian kernel is the predicted box position c_i . We compute the confidence weights over all positive predicted boxes. By incorporating the prior knowledge of layout from ground-truth, we can learn the weight for each predicted box. In Eq. 4, it illustrates that the spatial pattern score for predicted position c_i is a summation of weights by the ground truth positions which are inside the S^{c_i} . We compute the score over the input triples (c_i, N_i^*, u_i^*) :

$$K(c_i, N_i^*, u_i^*) = \begin{cases} \sum_{\theta \in D} \sum_{j \in N_i^*}^m G(j; \theta) & \text{if } u_i^* = 1\\ 1 & \text{otherwise,} \end{cases}$$
(4)

in which

$$G(j;\theta) = \alpha \cdot e^{-\left(\frac{x_j^{\theta}}{2\sigma_x^2} + \frac{y_j^{\theta}}{2\sigma_y^2}\right)},$$
(5)

is the 2D Gaussian spatial kernel that takes different rotated radius $D = \{\theta_1, ..., \theta_r\}$, where we use r = 4 ranged from 0 to π . The coordinate tuple (x_j, y_j) is the center position of *jth* ground-truth box in Eq. 5, and the coefficient α is the amplitude of the Gaussian function. All experiments use $\alpha = 1$.

We only give the weights for the foreground predicted box c_i where it is marked $u_i^* = 1$. By the means of aggregating weights from ground-truth boxes N_i^* in different direction kernels (Eq. 4), we can compute a summation of scores for taking various layout structures into account. It will give a higher probability to the object position, which has larger weight. Namely, the more similar objects of instances surrounding it, the more possible the predicted boxes are the same category of instances. Therefore, the predicted box collects the confidence from the same objects which are nearby (Figure 3). By leveraging spatially regularized weights, we can learn a model for generating the region proposals where the objects of instance will appear with their own layout.

5. Experiment

In this section, we evaluate our approach on two different datasets. The PUCPR+ dataset, made from the sub-dataset of the public PKLot dataset [10], and the CARPK dataset are both used to estimate the validation of our proposed LPNs. Then, we evaluate our object counting model, which leverages the structure information on the PUCPR+ dataset and our CARPK dataset.

5.1. Experiment Setup

We implement our model on Caffe [13]. For fairness in analyzing the effectiveness between different baseline methods, we implemented all of them based on the VGG-16 networks [25] which contains 13 convolutional layers and 3 fully-connected layers. All the layer parameters of baselines and our proposed model are using the weights pretrained on ImageNet 2012 [24], followed by fine-tuning the models on our CARPK dataset or the PUCPR+ dataset depending on the experiments. We run our experiments on the environment of Linux workstation with Intel Xeon E5-2650 v3 2.3 GHz CPU, 128 GB memory, and one NVIDIA Tesla K80 GPU. Our multi-task joint training scheme takes approximately one day to converge.

5.2. Evaluation of Region Proposal Methods

For evaluating the performance of our method LPNs, we use five-fold cross-validation on the PUCPR+ dataset to ensure that the same image would not appear across both training set and testing set. In order to better evaluate the recall while estimating localization accuracy, rather than reporting recall at particular IoU thresholds, we report the Average Recall (AR). It is an average of recall with IoU threshold t between 0.5 to 1, where $AR = \frac{1}{t} \sum_{i}^{t} Recall(IoU_t)$. As the metric of recall at IoU of 0.5 is not predictive of detection accuracy, proposals with high recall but at low overlap are not effective for detection [12]. Therefore, adopting the IoU range of the AR metric can simultaneously measure both proposal recall and localization accuracy to better predict the result of counting and localizing performance.

Table 2. Result on the PUCPR+ [10] dataset for average recall at 100, 300, 500, 700, and 1000 proposals with the different components of approaches. The method in the middle column represents the RPN training with the small default box size on conv4-3 layer.

#Proposals	RPN [22]	RPN+small	LPN (ours)
100	3.2%	20.5%	23.1%
300	9.1%	43.2%	49.3%
500	13.9%	53.4%	57.9%
700	17.4%	57.3%	60.7%
1000	21.2%	59.9%	62.5%

We compare our proposal method LPNs against the state-of-the-art object proposal generator RPNs [22] on the PUCPR+ dataset with different number of the object proposals. Our results are shown in Table 2. It reveals that our proposal method LPNs, which leverages the regularized layout information, can achieve higher recall and surpass RPNs in the same number of proposals. The state-ofthe-art object proposal RPNs suffer from poor performance in average recall. We refer that the factors, which affect the performance, are upon on the inappropriate anchor size and the resolution of convolutional layer features. Hence, in the same manner, we apply the smaller anchor box size on RPNs on the conv4-3 layer, which is in the same setting as our approach. Table 2 shows that the RPNs utilize the small anchor size and the higher resolution feature map bring about a better improvement. It implies that the CNN model is not as powerful in scale variance as we thought when using inappropriate layers or unsuitable default box size for prediction. This experiment also shows that the performance of our proposed model LPNs with spatial regularized constraints still outperforms the revised RPNs (e.g., 14.1% better in 300 proposals and 8.42% better in 500 proposals). Besides, we also found that our method with spatial regularizer significantly performs better in the dense case 2 . The result indicates that the prior layout knowledge could

 $^{^{2}}$ We additionally divide the PUCPR+ dataset into dense and less dense cases. Our method has 16.30% large relative improvement compared to RPN-small in dense case, which is better than 8.27% for less dense case. Moreover, our method localizes the bounding box more precisely, i.e., our method achieves 64.4% recall in IoU at 0.7 which is almost 10% better than RPN-small 54.7% for 300 proposals.

potentially benefit the outcome by giving the correct confidence score to the position of instances in images.

#Propo	osals R	PN [22]	RPN+small	LPN (ours)
100)	11.4%	31.1%	34.7%
300)	27.9%	46.5%	51.2%
500)	34.3%	50.0%	54.5%
700)	37.4%	51.8%	56.1%
100	0	39.2%	53.4%	57.5%

Table 3. Results on the CARPK dataset with different components.

For looking into the details of the effectiveness of our approach in region proposal, we also conduct the experiment on our CARPK dataset. In order to ensure that the same or the similar image scenes would not appear across both training and testing set, which would affect the observation of validation, we take 3 different scenes of the parking lot as training set and the remaining one scene of the parking lot as testing set. Table 3 reports the average recall of our methods, the state-of-the-art region proposal method RPNs, and the revised RPNs on CARPK dataset. In the experiment results, it comes as no surprise that by incorporating the additional layout prior information, our LPNs model boots both recall and localization accuracy of proposal method. Again, this result shows that the proposals generated by our approach are more effective and reliable.

5.3. Evaluation of Car Counting Accuracy

Since the goal of our proposed approach is to count the number of cars from drone view scenes, we compare our car counting system with three state-of-the-art methods which can also achieve the car counting task. A one-look regression-based counting method [18] which is the up-todate method for estimating the number of cars in density object counting measure and two prominent object detection systems, Faster R-CNN [22] and YOLO [21], which have remarkable success in object detection task in recent years.

For fair comparison, all the methods are built based on a VGG-16 network [25]. The only difference is that [18] uses a softmax with 64 outputs as they assumed that the maximum number of cars in a scene is sufficiently small. However, the maximum number of cars in the CARPK dataset is far more than 64. The maximum number of cars is 331 in a single scene of the PUCPR+ dataset and 188 in a single scene of the CARPK dataset. Hence, we follow the setting from [18] and train the network with a different output number to make this regression-based method compatible with the two datasets. We set the softmax to 400 outputs for the PUCPR+ dataset for evaluation. Last, the setting of two datasets, PUCPR+ and CARPK, are the same as the experiment of region proposal phase.

We employ two metrics, Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE), for evaluating the performance of counting methods. These two metrics have the similar physical meaning that estimates the error between the ground-truth car numbers y_i and the predicted car numbers f_i . MAE is the average absolute difference between ground-truth quantity y_i and predicted quantity f_i over all testing scenes where $MAE = \frac{1}{n} \sum_{i=1}^{n} |f_i| - \frac{1}{n} \sum_{i=1}^{n} |f_i|$ y_i . Similar, RMSE is the square root of the average of squared differences between ground-truth quantity and predicted quantity over all testing scenes where RMSE = $\sqrt{\frac{1}{n}\sum_{i}^{n}(f_{i}-y_{i})^{2}}$. The difference of the two metrics is that the RMSE should be more useful when large errors are particularly undesirable. Since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors. In the counting task, these metrics have good physical meaning for representing the average counting error of cars in the scene.

Table 4. Comparison with the object detection methods and the global regression method for car counting on the PUCPR+ dataset. N_p is the number of candidate boxes used in the object detector, which parameterizes the region proposal method. The "*" in front of the baseline methods represents that the method has been fine-tuned on PUCPR+ dataset. The "†" represents that the method is revised to fit our dataset.

Method	N_p	MAE	RMSE
YOLO [21]	-	156.72	200.54
Faster R-CNN [22]	400	156.76	200.59
*YOLO	-	156.00	200.42
*Faster R-CNN	400	111.40	149.35
*Faster R-CNN (RPN-small)	400	39.88	47.67
[†] One-Look Regression [18]	-	21.88	36.73
Our Car Counting CNN Model	400	22.76	34.46

We compare three methods on the PUCPR+ dataset, where the maximum number of cars is 311 in a single scene. Since the softmax output number of [18] is designed to be 400 for the PUCPR+ dataset, we also impartially compare to this dense object counting method with the number of region proposals limited to 400, which is a strict condition to our object counter. For YOLO, we select the parameter of confidence threshold at 0.15, which gives the best performance in our dataset.

The experimental results are shown in Table 4. The asterisk "*" in front of the YOLO and Faster R-CNN methods represents that the models have been fine-tuned on the PUCPR+ dataset, otherwise they are fine-tuned on the benchmark datasets (PASCAL VOC dataset and MS COCO dataset respectively), where they also have the car categories. Our proposed method outperforms the best RMSE on large-scale car counting, even with a very tough setting in the number of proposals. Note that we have comparable



Counting number: 292 cars Ground Truth: 299 cars

Counting number: 114 cars Ground Truth: 121 cars

Figure 5. Selected examples of car counting and localizing results on the PUCPR+ dataset (left) and the CARPK dataset (right). The counting model uses our proposed LPN trained on a VGG-16 model and combined with an object detector (Fast R-CNN), where the parameters setting of confidence score is 0.5 and non maximum suppression (NMS) is 0.3 for 2000 proposals.

MAE performance to the state-of-the-art car counting regression method [18], but the better RMSE implies that our method has better capability in some extreme cases. The methods that are fine-tuned on PASCAL and MS COCO get worse results. It reveals that the inferior outcomes are caused by the different perspective view of the object even when training with car category samples. The experiment results show that by incorporating the spatially regularized information, our Car Counting CNN model boosts the performance of counting. A counting and localization example result is shown in Figure 5 (left).

We further compare the counting methods on our challenging large-scale CARPK dataset where the maximum number of cars is 188 in a single scene. However, different from the PUCPR+ dataset which only has one parking lot, our CARPK dataset provides various scenes of diverse parking lots for cross-scene evaluation. In the setting of [18] method, we also deign a 200 softmax output network for the CARPK dataset. In order to fairly compare the counting methods, we again restrict the number of proposals of object counter which has utilized the region proposal method with a tough number 200^{-3} . The quantitative results of car counting on our dataset are reported in Table 5. The experiment results show that our car counting approach is reliably effective and has the best MAE and RMSE even in the cross-scene estimation. An counting and localizing example result is shown in Figure 5 (right). Still our method can generate the feasible proposals and obtain the reasonable counting result close to the real number of cars in the

scenes of the parking lots.

Table 5. Comparison results on the CARPK dataset. The notation definition is similar to Table 4.

Method	N_p	MAE	RMSE
YOLO [21]	-	102.89	110.02
Faster R-CNN [22]	200	103.48	110.64
*YOLO	-	48.89	57.55
*Faster R-CNN	200	47.45	57.39
*Faster R-CNN (RPN-small)	200	24.32	37.62
†One-Look Regression [18]	-	59.46	66.84
Our Car Counting CNN Model	200	23.80	36.79

6. Conclusions

We have created the to-date largest drone view dataset, called CARPK. It is a challenging dataset for various scenes of parking lots in a large-scale car counting task. Also, in the paper, we introduced a new way for generating the feasible region proposals, which leverage the spatial layout information for an object counting task with regularized structures. The learned deep model can specifically count objects better with the prior knowledge of object layout patterns. Our future work will involve global information, such as context, road scene, and other objects which can help distinguish between false car-like instances and real cars.

7. Acknowledgement

This work was supported in part by the Ministry of Science and Technology, Taiwan, under Grant MOST 104-2622-8-002-002 and MOST 105-2218-E-002-032, and in part by MediaTek Inc, and grants from NVIDIA and the NVIDIA DGX-1 AI Supercomputer.

 $^{^{3}}$ Our method gets better performance when using bigger number of proposals (e.g., 8.04 and 12.06 for 1000 proposals in MAE and RMSE respectively) in the PUCPR+ dataset. In the CARPK dataset, our method also has 13.72 and 21.77 for 1000 proposals in MAE and RMSE respectively.

References

- M. Ahrnbom, K. Astrom, and M. Nilsson. Fast classification of empty and occupied parking spaces using integral channel features. In *CVPR*, 2016. 1
- [2] G. Amato, F. Carrara, F. Falchi, C. Gennaro, and C. Vairo. Car parking occupancy detection using smart camera networks and deep learning. In *ISCC*, 2016. 1
- [3] S. An, W. Liu, and S. Venkatesh. Face recognition using kernel ridge regression. In CVPR, 2007. 3
- [4] P. Arbeláez, J. Pont-Tuset, J. T. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping. In CVPR, 2014. 3
- [5] C. Arteta, V. Lempitsky, J. A. Noble, and A. Zisserman. Interactive object counting. In *ECCV*, 2014. 3
- [6] A. B. Chan, Z.-S. J. Liang, and N. Vasconcelos. Privacy preserving crowd monitoring: Counting people without people models or tracking. In *CVPR*, 2008. 3
- [7] K. Chen, S. Gong, T. Xiang, and C. Change Loy. Cumulative attribute space for age and crowd density estimation. In *CVPR*, 2013. 3
- [8] K. Chen, C. C. Loy, S. Gong, and T. Xiang. Feature mining for localised crowd counting. In *BMVC*, 2012. 3
- [9] M.-M. Cheng, Z. Zhang, W.-Y. Lin, and P. Torr. Bing: Binarized normed gradients for objectness estimation at 300fps. In *CVPR*, 2014. 3
- [10] P. R. de Almeida, L. S. Oliveira, A. S. Britto, E. J. Silva, and A. L. Koerich. Pklot–a robust dataset for parking lot classification. *Expert Syst Appl*, 2015. 1, 2, 3, 6
- [11] R. Girshick. Fast r-cnn. In ICCV, 2015. 5
- [12] J. Hosang, R. Benenson, P. Dollár, and B. Schiele. What makes for effective detection proposals? *TPAMI*, 2016. 6
- [13] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 6
- [14] D. Kamenetsky and J. Sherrah. Aerial car detection and urban understanding. In *DICTA*, 2015. 3
- [15] D. Kong, D. Gray, and H. Tao. A viewpoint invariant approach for crowd counting. In *ICPR*, 2006. 3
- [16] V. Lempitsky and A. Zisserman. Learning to count objects in images. In NIPS, 2010. 3
- [17] T. Moranduzzo and F. Melgani. Automatic car counting method for unmanned aerial vehicle images. *IEEE Transactions on Geoscience and Remote Sensing*, 2014. 3
- [18] T. N. Mundhenk, G. Konjevod, W. A. Sakla, and K. Boakye. A large contextual dataset for classification, detection and counting of cars with deep learning. In *ECCV*, 2016. 1, 2, 3, 7, 8
- [19] P. O. Pinheiro, R. Collobert, and P. Dollar. Learning to segment object candidates. In *NIPS*, 2015. 3
- [20] S. Razakarivony and F. Jurie. Vehicle detection in aerial imagery: a small target detection benchmark. *JVCIR*, 2016. 2, 3
- [21] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 7, 8

- [22] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 2, 3, 4, 5, 6, 7, 8
- [23] M. Rodriguez, I. Laptev, J. Sivic, and J.-Y. Audibert. Density-aware person detection and tracking in crowds. In *ICCV*, 2011. 3
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 6
- [25] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556, 2014. 6, 7
- [26] V. A. Sindagi and V. M. Patel. Generating high-quality crowd density maps using contextual pyramid cnns. In *ICCV*, 2017.
 3
- [27] C. Szegedy, S. Reed, D. Erhan, and D. Anguelov. Scalable, high-quality object detection. arXiv preprint arXiv:1412.1441, 2014. 3
- [28] F. Tanner, B. Colder, C. Pullen, D. Heagy, M. Eppolito, V. Carlan, C. Oertel, and P. Sallee. Overhead imagery research data setan annotated data library & tools to aid in the development of computer vision algorithms. In *AIPR*, 2009. 2, 3
- [29] J. R. Uijlings, K. E. van de Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *IJCV*, 2013. 3
- [30] M. Wang and X. Wang. Automatic adaptation of a generic pedestrian detector to a specific traffic scene. In *CVPR*, 2011.
- [31] C. Zhang, H. Li, X. Wang, and X. Yang. Cross-scene crowd counting via deep convolutional neural networks. In *CVPR*, 2015. 3
- [32] S. Zhang, G. Wu, J. P. Costeira, and J. M. Moura. Fcn-rlstm: Deep spatio-temporal neural networks for vehicle counting in city cameras. In *ICCV*, 2017. 3
- [33] C. L. Zitnick and P. Dollár. Edge boxes: Locating object proposals from edges. In ECCV, 2014. 3