# A Self-Training Approach for Point-Supervised Object Detection and Counting in Crowds

Yi Wang, *Graduate Student Member, IEEE*, Junhui Hou, *Senior Member, IEEE*, Xinyu Hou, *Student Member, IEEE*, and Lap-Pui Chau, *Fellow, IEEE*

*Abstract*—In this paper, we propose a novel self-training approach named Crowd-SDNet that enables a typical object detector trained only with point-level annotations (i.e., objects are labeled with points) to estimate both the center points and sizes of crowded objects. Specifically, during training, we utilize the available point annotations to supervise the estimation of the center points of objects directly. Based on a locally-uniform distribution assumption, we initialize pseudo object sizes from the point-level supervisory information, which are then leveraged to guide the regression of object sizes via a crowdedness-aware loss. Meanwhile, we propose a confidence and order-aware refinement scheme to continuously refine the initial pseudo object sizes such that the ability of the detector is increasingly boosted to detect and count objects in crowds simultaneously. Moreover, to address extremely crowded scenes, we propose an effective decoding method to improve the detector's representation ability. Experimental results on the WiderFace benchmark show that our approach significantly outperforms state-of-the-art point-supervised methods under both detection and counting tasks, i.e., our method improves the average precision by more than 10% and reduces the counting error by 31.2%. Besides, our method obtains the best results on the crowd counting and localization datasets (i.e., ShanghaiTech and NWPU-Crowd) and vehicle counting datasets (i.e., CARPK and PUCPR+) compared with state-of-the-art counting-by-detection methods. The code will be publicly available at https://github.com/WangyiNTU/Point-supervised-crowd-detection.

*Index Terms*—Convolutional neural network (CNN), object detection, crowd counting, self-training, weak supervision.

## I. INTRODUCTION

With a massive amount of population living in cities, crowd scenes have become a fundamental yet challenging scenario in a wide variety of computer vision applications, such as video surveillance [1], crowd analysis [2], [3], and safety monitoring [4], [5]. Objects in dense crowds present small sizes, large scale variations, and high occlusions, which poses great challenges to object detection methods that simultaneously predict objects' locations and sizes in an image.

The advances of deep neural networks (DNNs) raise an issue of enormous demand for data annotations. However, it is very costly and laborious to collect object-level bounding box annotations [6], [7] which are usually needed for training DNN-based object detection methods, especially for

Yi Wang, Xinyu Hou, and Lap-Pui Chau are with School of Electrical and Electronics Engineering, Nanyang Technological University, Singapore, 639798 (e-mail: {wang1241, houx0008}@e.ntu.edu.sg, elpchau@ntu.edu.sg).

Junhui Hou is with the Department of Computer Science, City University of Hong Kong (e-mail: jh.hou@cityu.edu.hk).
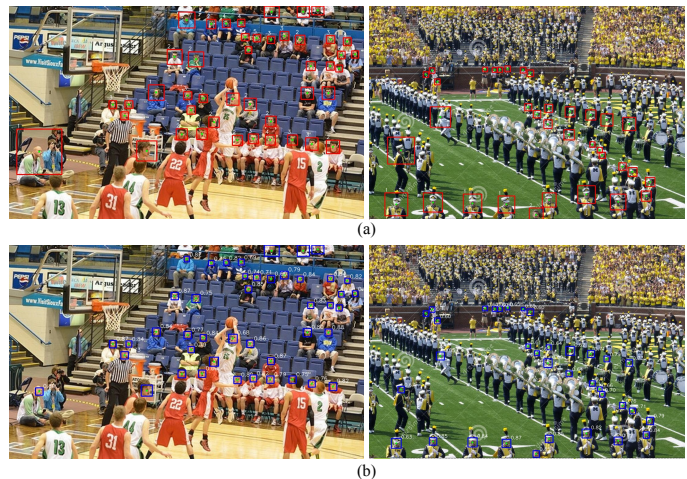
Corresponding author: Lap-Pui Chau.



Fig. 1. Illustration of generated training examples at different phases by our self-training approach. The images are from WiderFace dataset [12]. The green dots shown at the centers of faces stand for the point-level annotations. (a): Before training, the pseudo object sizes (red boxes) generated by the proposed locally-uniform distribution assumption. The numbers shown at the top-right corner of red boxes stand for object crowdedness. (b): After training, the refined pseudo object sizes (blue boxes) by our crowdedness-aware loss and the confidence and order-aware refinement scheme. The numbers shown at the top-left corner of blue boxes stand for the pseudo object sizes' posterior probabilities. *Zoom in the figure for better viewing.*

images containing thousands of objects. Current crowd counting datasets provide only point-level annotations, and usually human heads are labeled as the central points, e.g., the green dots shown in Fig. 1. Due to the lack of object sizes, state-of-the-art DNN-based object detectors [8] cannot be trivially applied to such point supervision. As pioneers, Liu *et al.* [9] introduced a pseudo size updating scheme in a detection network to estimate object sizes. Sam *et al.* [10] proposed an LSC-CNN to achieve higher detection performance in crowd scenes. However, these works are still not on par with box-supervised methods (e.g., Faster R-CNN [11]) in the detection task. As for the counting task, these methods, denoted as counting-by-detection methods, can count objects by filtering out low-confidence objects with a threshold. However, they also suffer from highly crowded objects.

Alternatively, modern crowd counting methods [13]–[15], named counting-by-regression methods, bypass the locations and sizes of objects but employ DNNs to regress a density map, which is further integrated to obtain the overall count. These methods have achieved outstanding counting performance in dense crowds. However, they only aim to count

the number of objects and lose individual information, i.e., object instance's location and size. We argue that the density map only contains weak information about the crowds, while locations and sizes of object instances provide much more important information for other computer vision applications, such as multi-object tracking [16], [17], face recognition [12], [18], and person re-identification [19].

In view of these issues, we propose a novel self-training approach capable of training a typical detection method only with point-level annotations such that it can accurately and simultaneously detect and count objects in dense crowds. Specifically, based on a keypoint-based detector, i.e., center and scale prediction (CSP) [20], we decouple detection as the separate estimation of objects' center points and sizes. The available point-level annotations directly supervise the estimation of the center points during training. As the ground-truth object sizes are not accessible, we propose a simple yet effective assumption in crowd scenes, called locally-uniform distribution assumption (LUDA), to generate the initial pseudo size for each object (see the red bounding boxes shown in Fig. 1(a)). Meanwhile, we propose a crowdedness-aware loss to emphasize the contributions of crowded objects in object size regression (see the crowdedness at the top-left corner of red boxes in Fig. 1(a)). Moreover, we propose a confidence and order-aware refinement scheme to continuously update the pseudo sizes during training, which performs the refinement operation by considering both the prior confidences and the updating order of pseudo sizes, such that the detection ability of the detector is increasingly boosted (see the blue bounding boxes shown in Fig. 1(b) for the refined pseudo object sizes). Besides, to deal with highly dense crowds (e.g., one person represented by several pixels in an image on the ShanghaiTech [21] dataset), we propose an effective decoding method to improve the representation ability of the detector, in which a feature fusion and decoding technique is employed to restore the full-resolution feature maps.

Extensive experimental results show that our approach outperform start-of-the-art point-supervised methods to a significant extent in terms of the detection performance, i.e., more than 10% AP improvement is achieved. Moreover, our method even produces comparable performance to the box-supervised Faster R-CNN on the WiderFace benchmark [12]. For center point localization, our approach produces the best results among state-of-the-art methods on dense crowd datasets, e.g., ShanghaiTech [21] and NWPU-Crowd [22]. For crowd counting, our method obtains comparable results to the latest counting-by-regression methods. Note that the bounding boxes produced by our method are more informative than the density map produced by counting-by-regression methods.

The rest of this paper is organized as follows. In Section II, we introduce the related works on object detection and counting in crowds. Then, the proposed self-training approach is presented in detail in Section III. Experimental results and ablation studies are provided in Section IV. Finally, we conclude this paper in Section V.

## II. RELATED WORK

### A. Object Detection in Dense Crowds

Benefiting from the advances of DNNs, recent object detectors, such as Faster R-CNN [11], RetinaNet [23], and CenterNet [24], [25], achieve appealing performance. Despite the remarkable progress made, these methods encounter difficulties when counting small and heavily occluded objects in crowded scenes. To enhance the detectors' abilities, Liu *et al.* [20] proposed a keypoint-based detector named CSP to predict the central points and scale of objects separately. Goldman *et al.* [26] presented a deep-learning-based method for precise object detection in densely packed scenes. They introduced a soft intersection-over-union (IoU) network layer and an expectation-maximization (EM) based clustering method to deal with overlapped objects in the dense scenes. Though these above detectors all achieve good performance, they have to be trained and supervised by box-annotated examples.

Most crowd counting datasets only provide point annotations for denser crowds. It is nontrivial to train a detector with point supervision. Similarly, sample weighting techniques [27] that require the box supervision fail to work on point-supervised detection. Recently, several works begin to use the blobs [28] to localize the individuals in crowds and even to estimate the sizes of the heads [9], [10], with only point-level annotations. Laradji *et al.* [28] argued that the unnecessary size and shape information drags the performance of detection-based methods in counting problems. A localization-based counting loss that combines image-level, point-level, split-level, and false-positive loss is used to train a fully convolutional network (FCN) such that it could produce the blobs in the center of objects. Based on a regression-based network, Idrees *at al.* [29] proposed a post-processing method to find the local peaks on the density map as the center locations of heads. As a baseline method, Sam *et al.* [10] applied a threshold technique on the density map of the CSRNet [13] to obtain detections (called CSR-A-thr). However, these methods only localize the center points of individuals in crowds.

To further estimate individuals' sizes, Liu *et al.* [9] proposed a detection network, named PSDDN, which builds a strong baseline for point-supervised detection and counting in crowds. The PSDDN employs the nearest neighbor distance [21] to initialize the pseudo boxes and updates the pseudo boxes by choosing smaller box predictions. Another state-of-the-art method, named LSC-CNN, was proposed recently by Sam *et al.* [10], '

### B. Object Counting in Dense Crowds

Bypassing localization, counting-by-regression methods [1], [21], [30], [31] were proposed to address the crowd counting problem and have dominated this field for years. Instead of directly regressing the global count adopted in early works [4], [32], current approaches [13], [14], [33] exploit DNNs to estimate a density map [34], over which the count is obtained via the integration operation. The typical network architecture of this kind of methods is generally composed of an encoder for extracting a set of features from an image and a decoder for regressing a density map [35]. For example, in [21], [30],
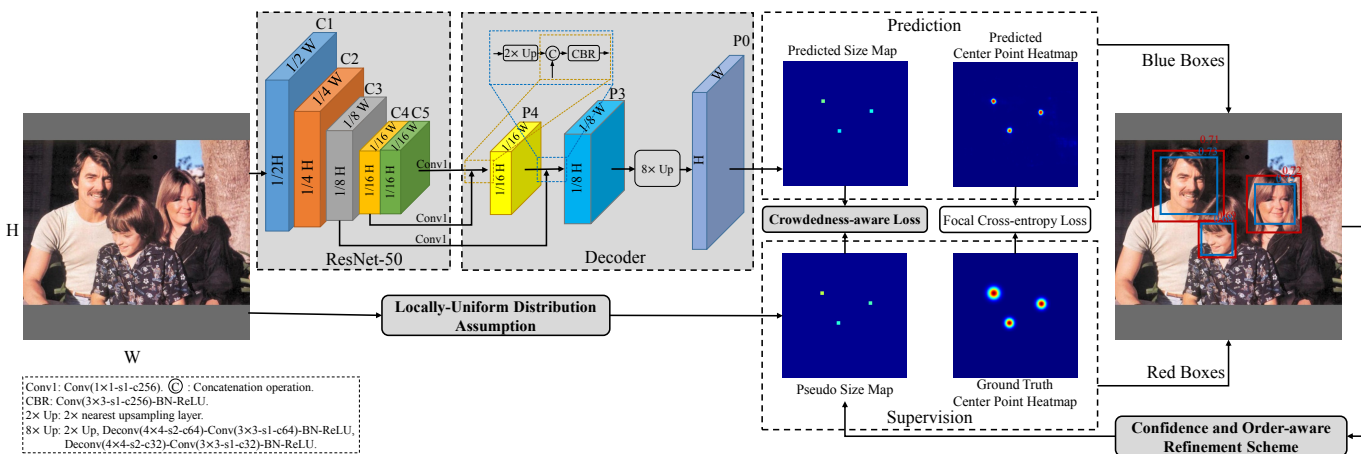
Fig. 2. Illustration of the proposed self-training framework for training a detector only with point-wise supervision. The decoder produces high-resolution feature maps for the estimation of objects' center points and sizes. Before training, the pseudo object sizes are first generated from the point-wise supervisory information, based on our locally-uniform distribution assumption. During training, the pseudo object sizes are further refined by our confidence and order-aware refinement scheme under the supervision of the proposed crowdedness-aware loss. Here an image with sparse objects is used for clear visualization purposes.

a multi-column structure was introduced to learn the multi-scale features for representing large scale variation of objects. Li *et al.* [13] demonstrated the redundancy of features in the multi-column structure and proposed a single-column deep structure (e.g., VGG [36]) with dilated convolution layers, which achieved better results. Cao *et al.* [33] proposed a scale aggregation module to learn the scale diversity of features. Shi *et al.* [14] repurposed point annotations as a segmentation map and a global density map.

By exploiting the attention mechanism, Sindagi *et. al.* [37] proposed an attention-based crowd counting network, named HA-CCN, to enhance the different-level features of the network. To realize a lightweight model, Wang *et al.* [35] proposed a pixel shuffle decoder (PSD) to generate a high-resolution density map without using convolutional layers. They presented a density-aware curriculum learning (DCL) strategy to improve the network's generalization ability and reduce training time. Recently, Wang *et al.* [22] provided a large-scale benchmark, namely NWPU-Crowd, for the crowd localization and counting. A recent survey paper written by Gao *et al.* [15] reviews over two-hundred crowd counting works and some datasets, showing the improvement in this filed.

Although the counting-by-regression methods obtain state-of-the-art counting performance, they sacrifice the location and size information of objects [15]. We argue that the count is rough information for the crowds and not enough for further use of high-level vision tasks. In our work, we count objects by our object detections.

### C. Counting from Drone View

Another dense case appears in vehicle counting from drone images. Hsieh *et al.* [38] introduced a large-scale car parking lot dataset (CAPPK), which consists of approximately 90k cars. Motivated by the regular spatial layout of cars, they proposed a layout proposal network (LPN) to count vehicles. Without local annotations such as bounding boxes or points,

Stahl *et al.* [39] proposed a general object counting method that uses local image regions to predict the global image-level counts. Goldman et al. [26] considered that vehicles in drone images are a densely packed scene. Therefore, they used their object detection method in this scenario. To deal with the failure of state-of-the-art detectors in drone scenes, Li *et al.* [40] made a set of modifications for detection. An effective loss is proposed to yield the scale-adaptive anchors. Then, the circular flow is applied to guide feature extraction. Third, a counting regularized constraint is introduced to the loss function.

## III. PROPOSED METHOD

Fig. 2 illustrates the proposed self-training framework, which is capable of training a typical object detector only with point-level annotations for simultaneous object detection and counting. To be specific, the framework is based on an anchor-free keypoint-based object detector, i.e., CSP detector [20]. A locally-uniform distribution assumption is proposed to generate the initial pseudo size for each object, and a crowdedness-aware loss is proposed to emphasize the pseudo sizes of crowded objects in size regression. Furthermore, a confidence and order-aware refinement scheme is proposed to update the pseudo sizes in each training iteration. In addition, a decoding method is proposed to handle dense crowds. In what follows, we present the detection network and the self-training approach in detail.

### A. Detection Network

We employ the keypoint-based detection method because it allows us to estimate the center point and size of objects separately. Therefore, the center points could be directly supervised by the point-level annotation, while the size estimation is accomplished by the proposed modules. In addition, it is anchor box-free.

*1) Architecture:* For the datasets with medium-density crowds (e.g., WideFace [12], CARPK and PUCPR+ [38]), we employ the original network in [20]. For the dataset with high-density crowds (e.g., ShanghaiTech [21] and NWPU-Crowd [22]), however, the original network performs poorly since the center point prediction fails to represent highly dense objects in its output feature map of size $\frac{H}{4} \times \frac{W}{4}$, where $H$ and $W$ are the height and width of the input image, respectively. Hence, we propose an effective decoding method to handle this issue.

As illustrated in Fig. 2, the detection network contained in our approach adopts the five-stage ResNet-50 [41] as the backbone network, where each stage downsamples the feature maps by a factor of 2, except for Stage-5 that uses the dilated convolutions to keep the stride the same as Stage-4. Let $C_i, i \in \{1, 2, ..., 5\}$ be the output feature maps of the $i$-th stage. In the *Decoder*, three Conv(1×1-s1-c256) are applied to reduce the number of channels of $C_5$, $C_4$, and $C_3$, where Conv denotes the convolutional layer, and (1×1-s1-c256) means the layer with the kernel of 1×1, the stride of 1, and the channel of 256. Then, we employ a top-down feature fusion manner to merge $C_5$, $C_4$, and $C_3$, generating the fused features $P_3$. The fusion manners are shown as the dotted yellow and blue rectangle in Fig. 2. Finally, we use an 8× Up structure to decode the fused features, consisting a 2× nearest upsampling layer followed by two Deconv(4×4-s2)-Conv(3×3-s1)-BN-ReLU, where Deconv, BN, and ReLU denote the deconvolution, batch normalization, rectified linear unit, respectively. The decoder produces the output feature maps ($P_0$) with the same size as the input, i.e., $H \times W$. There are two separate heads (i.e., Conv(1×1-s1-c1)) for center point prediction and size prediction, producing the center point heatmap and the size map.

*2) Supervision information:* Let $\{\mathbf{p}_j\}_{j=1}^{M}$ be the point-wise annotations, where $p_j := (x_j, y_j)$ is the 2D coordinates of the center of the $j$-th object in an image, and $M$ is the total number of objects in an image. As shown in Fig. 2, to supervise the estimation of object center points, we generate a ground-truth center point heatmap $Q \in [0, 1]^{H \times W}$ with 1 for objects' center points and 0 for negative points. To decrease the ambiguity of negative points surrounding the positive ones, we place a normalized 2D Gaussian mask at the center location of each positive point, as performed in [20]. If two masks overlap, we choose the element-wise maximum for the overlapped region. For the object size supervision, we generate pseudo object sizes denoted by $s_j$, which will be introduced in Sec. III-B. Here we assume that the objects (e.g., heads and faces) have an aspect ratio of 1 in crowded scenes. We assign $log(s_j)$ to the $j$-th object's center coordinates $(x_j, y_j)$ and zeros to other locations, generating the size map $S \in \mathbb{R}^{H \times W}$. For the original CSP with the output of size $\frac{H}{4} \times \frac{W}{4}$, an offset map is appended to estimate the discretization error caused by the stride of 4. The ground-truth offset for $p_j$ is defined as $\frac{x_j}{4} - \lfloor \frac{x_j}{4} \rfloor$ and $\frac{y_j}{4} - \lfloor \frac{y_j}{4} \rfloor$ on the $x$-axis and $y$-axis, respectively, which is assigned to $(x_j, y_j)$ on the offset map. $\lfloor r \rfloor : \mathbb{R} \to \mathbb{Z}$ of a real number $r$ denotes the greatest integer less than or equal to $r$.

*3) Loss function:* We apply the focal cross-entropy loss [20], [24] to each pixel on the center point heatmap:

$$L_c = -\frac{1}{M} \sum_{j=1}^{M} \begin{cases} (1 - \hat{q}_j)^\gamma log(\hat{q}_j), & \text{if } q = 1, \\ A(1 - q_j)^\delta (\hat{q}_j)^\gamma log(1 - \hat{q}_j), & \text{otherwise,} \end{cases} \tag{1}$$

where $\hat{q}_j$ and $q_j$ are the predicted probability and the ground-truth label of pixel $j$, respectively; $\gamma$ is the hyper-parameter of the focal loss [23] which is set to 2 in all experiments; $\delta$ is the hyper-parameter to control the penalty of negatives, which is set to 4 in all experiments; $A$ is the coefficient to address the imbalance between positive and negative points, which is set to 1 (resp. 1/16) for the original CSP (resp. the proposed decoding structure) in all experiments. Intuitively, if we enlarge the CSP's output feature map by the decoding method from $\frac{H}{4} \times \frac{W}{4}$ to $H \times W$, the number of negative points will increase by 16 times. Hence, we set $A = 1/16$ to balance the positive and negative points. For object size regression, we propose a crowdedness-aware loss $L_{size-\alpha}$, of which the details will be described in Sec. III-C. For offset estimation, the smooth L1 loss [42] is calculated between the ground-truth offsets and predicted ones, denoted as $L_o$. The overall training objective is

$$L = \lambda L_c + L_{size-\alpha} + L_o, \tag{2}$$

where $\lambda$ is the weight for center point classification, which is experimentally set to 0.1.

*4) Inference:* The detector first performs a forward pass to generate the center point heatmap, the size map, and the offset map (if it is used). The peak points in the center point heatmap are extracted by a 3×3 max pooling operation. Then, we obtain the center point coordinates of objects whose probabilities are larger than a predefined confidence. The object sizes are obtained from the corresponding coordinates in the size map (see the "Prediction" in Fig. 2). If the offset map is appended, the corresponding offsets are added to object coordinates. Finally, the bounding boxes can be decoded by the coordinates and sizes.

### B. LUDA-based Pseudo Object Size Generation

As the ground-truth object sizes are unavailable, we generate pseudo object sizes from the point-wise supervisory information to train the detector. In crowded scenarios, object instances, e.g., heads, faces, or cars, are usually uniformly distributed in an image. According to the geometry-adaptive kernel [21], a typical object's size is proportional to the distance to its nearest neighbors in dense crowds. This assumption is relatively weak as the objects are not always dense enough in an image. In this paper, we employ the non-uniform kernel [14] to locally restrict the above assumption and propose the locally-uniform distribution assumption (LUDA). That is, 1) the objects in crowd scenarios are uniformly distributed in a local region and have a similar size in that region; and 2) the crowdedness of the region affects the precision of size estimates. In what follows, we detail our pseudo object size generation method based on LUDA.

Following [14], [21], we first calculate the initial object size of point $\mathbf{p}_j$ according to the distances to its $K$ nearest points, i.e.,

$$\overline{d}_j = \frac{1}{K} \sum_{k=1}^{K} \beta d_{j,k}, \qquad (3)$$

where $\overline{d}_j$ is the initial object size of $\mathbf{p}_j$, $d_{j,k}$ is the distance between point $\mathbf{p}_j$ and its $k$-th nearest neighbor, and $\beta$ is a scalar. The initial object size is further smoothed to reduce the variation in a local region, leading to the pseudo object size $s_j$:

$$s_j = \frac{1}{\left| R_{p_j} \right|} \sum_{l \in R_{p_j}} \overline{d}_l, \qquad (4)$$

where $R_{p_j}$ is the $p_j$-centered local region, $\left| R_{p_j} \right|$ is the number of the points inside $R_{p_j}$, and $\overline{d}_l$ is the initial size of the $l$-th point contained in $R_{p_j}$. We set the $R_{p_j}$ to a circular region so that KD-Tree can be used to improve the calculation speed. $\left| R_{p_j} \right|$ affects the precision of size estimates, which will be further exploited for the crowdedness-aware loss in Sec. III-C.

Since the crowdedness differs from regions to regions, a single set of parameters in Eqs. (3) and (4) can only fit a specific crowdedness. In experiments, we choose multiple sets of parameters intuitively for different training sets such that the crowded objects have precise pseudo sizes. The detailed parameter settings will be explained in Sec. IV-A2.

### C. Crowdedness-aware Loss for Object Size Regression

With the pseudo object size $s_j$, a straightforward way to supervise the regression of object sizes is using the smooth L1 loss [42], i.e.,

$$L_s = \frac{1}{M} \sum_{j=1}^{M} SmoothL1(\hat{s}_j, s_j), \qquad (5)$$

where $\widehat{s}_j$ is the size prediction of the $j$-th object. However, there is an obvious drawback for Eq. (5), i.e., all object instances equally contribute to the loss, and if some pseudo object sizes are inaccurate (i.e., noisy supervision), the training of the detector will be adversely affected. Moreover, inaccurate pseudo sizes are inevitable since the simple LUDA cannot resolve the complexity of crowded scenes (see Fig. 1(a)). Such a drawback is experimentally verified in the following Table III (see the 5th entry of Table III).

To this end, we propose a crowdedness-aware loss associating each object instance with its crowdedness, which indicates the importance in object size regression. We formulate the crowdedness-aware loss as

$$L_{size-\alpha} = \frac{1}{M} \sum_{j=1}^{M} \alpha_j SmoothL1(\widehat{s}_j - s_j), \qquad (6)$$

where $\alpha_j \in [0, +\infty)$ is the crowdedness-aware factor controlling weight of the $s_j$. Based on LUDA, crowded objects can produce more accurate pseudo object sizes than sparse ones, and thus the crowded objects should be assigned with larger weights. Specifically, we define the $\alpha_j$ as an exponential

function of the number of objects (crowdedness) inside the local region $R_{p_j}$, i.e.,

$$\alpha_j = (\left| R_{p_j} \right|)^{\eta}. \qquad (7)$$

Here, we use a tunable parameter $\eta \geq 0$ to scale the factor. In practice, we use a threshold to limit the maximum value of $\alpha_j$ to avoid gradient explosion. By assigning crowdedness-aware factors to the pseudo sizes, the loss function emphasizes the influence of the crowded objects and weakens that of the sparse ones. The crowdedness-aware loss enhances the robustness of the training with the noisy pseudo sizes.

### D. Confidence and Order-aware Refinement Scheme for Pseudo Size Updating

Although the crowdedness-aware loss is able to boost the robustness of the detector, the inaccurate sizes still exist and will affect the backward pass of training. Thus, we propose a confidence and order-aware refinement scheme to update pseudo object sizes for better training the detector. In [9], the pseudo bounding boxes are updated by selecting the predicted boxes with the highest scores among those whose sizes are smaller than the pseudo ones. However, such a criterion may not be true in practice since it ignores the following two key issues, resulting in the inaccurate refinement and unstable training process. 1) The prior information. It updates the pseudo sizes without considering their prior confidences. 2) The updating order. All pseudo sizes are treated identically, resulting in the synchronous updating of both accurate and inaccurate sizes.

Instead, by taking prior information into account, we assign every pseudo size with a prior probability at the beginning of training, and if and only if the predicted posterior probability of the detector is larger than the prior probability, we update the pseudo size (resp. prior probability) with the predicted size (resp. posterior probability) for the next epoch. Referring to the prior confidence, our refinement scheme guarantees the detector is trained with increasingly confident examples. Meanwhile, it can update the most inaccurate sizes first, then followed by updating the relatively accurate ones. This is achieved by setting the same prior probability for all pseudo sizes before training. During training, easy examples (e.g., sparse and large objects) with noisy size supervision that achieve high posterior probabilities rapidly are updated first. Then, hard examples (e.g., crowded and small objects) with more accurate size supervision are updated. In other words, the updating order is from noisy sizes to noiseless sizes. The proposed refinement scheme presents a more powerful error-correction capability than the one in [9]. See the experimental results in Table III.

More specifically, let $b_j = \{p_j, s_j\}$ be a pseudo bounding box centered at $p_j$ with the pseudo size $s_j$, and the prior probability $P(b_j)$ is assigned to $b_j$. The object detection problem is cast as learning the posterior $P(C, B|X)$, where $X$ is the input image, $B$ is the bounding boxes of objects, and $C \in \{0, 1\}$ is the binary class with 0 for background and 1 for object instance. Our refinement scheme, which is merged into the training process, is summarized in Algorithm 1. In each

**Algorithm 1** Confidence and Order-aware Refinement Scheme.

---

**Input:** The $i$-th input image $X_i$ with a set of pseudo bounding boxes $B_i = \{b_1, ..., b_{M_i}\}$, and prior probabilities of the boxes $P(b_j), j \in \{1, ..., M_i\}$.

**Output:** Refined pseudo bounding boxes $\widetilde{B}_i = \{\widetilde{b}_1, ..., \widetilde{b}_{M_i}\}$, and refined prior probabilities $P(\widetilde{b}_j)$.

---

1: Forward passing the detector to generate the detections $\widehat{B}_i = \{\widehat{b}_1, ..., \widehat{b}_{M_i}\}$ with the posterior probabilities $P(C, \widehat{B}_i | X_i)$.
2: **for** $j$ in $\{1, ..., M_i\}$ **do**
3:    **if** $P(C, \widehat{b}_j | X_i) > P(b_j)$ **then**
4:       $P(\widetilde{b}_j) \leftarrow P(C, \widehat{b}_j | X_i)$
5:       $\widetilde{b}_j \leftarrow \widehat{b}_j$
6:    **else**
7:       $P(\widetilde{b}_j) \leftarrow P(b_j)$
8:       $\widetilde{b}_j \leftarrow b_j$

---

training iteration, after executing a forward path, we obtain the predicted detections $\widehat{B}_i$ with the posterior $P(C, \widehat{B}_i | X_i)$ (i.e., Line 1). For each instance $j \in \{1, ..., M_i\}$, if $P(C, \widehat{b}_j | X_i)$ is larger than $P(b_j)$, we update the prior probabilities $P(b_j)$ with the posterior $P(C, \widehat{b}_j | X_i)$, and update the pseudo bounding boxes $\widetilde{b}_j$ with the prediction $\widehat{b}_j$ (i.e., Lines 3-5). Otherwise, the prior probability and pseudo bounding box remain unchanged (i.e. Lines 7 and 8). A constant prior probability of 0.6 is assigned to all $\{B_i\}_i^N$ at the beginning of training, where $N$ is the number of training images.

*Remark.* Here we provide more explanations why our self-training approach enables the detector to update the pseudo size towards the correct direction. The main reasons come from two aspects. First, as analyzed in Sec. III-B, our LUDA-based pseudo size generation method ensures that the pseudo sizes of crowded objects (i.e., hard examples) are more accurate than those of sparse objects. Also, it is known that hard examples have a heavy influence on the detector's training than easy examples [43], [44]. Thus, when incorporated with the crowdedness-aware loss, the larger number of relatively accurate hard examples in crowded scenes dominate the detector, making it robust to the outliers (e.g., inaccurate pseudo sizes). Second, the pseudo sizes are updated in a robust and orderly manner by our refinement scheme. We set the same initial prior probability for all pseudo bounding boxes. The posterior probability of easy examples will first meet the initial prior probability during training, and thus Algorithm 1 begins with updating the pseudo sizes of inaccurate easy examples. With increasingly accurate easy examples, the detector becomes stronger such that hard examples are then refined. More experimental results to illustrate the effectiveness of the proposed method can be found in Sec. IV-B.

## IV. EXPERIMENTAL RESULTS

In this section, we first describe the experiment settings, including datasets, implementation details, and evaluation metrics, followed by ablation studies for verifying the effectiveness of each component of our approach. Finally, we compare our approach with state-of-the-art methods in terms of both detection and counting tasks.

TABLE I
THE PARAMETER SETTINGS OF THE LUDA-BASED PSEUDO SIZE GENERATION METHOD FOR DIFFERENT DATASETS. "-" MEANS THE PARAMETER IS NOT REQUIRED.

| Dataset | Generation Parameters | | |
|---|---|---|---|
| | $K$ of Eq. (3) | $\beta$ of Eq. (3) | Max $\alpha_j$ of Eq. (7) |
| WiderFace [12] | 2 | 0.5 | - |
| SHA [21] | 2 | 1 | 50 |
| SHA [21] | 2 | 0.5 | 50 |
| NWPU-Crowd [22] | 2 | 0.8 | 200 |
| CARPK [38] | 1 | 1.3 | - |
| PUCPR+ [38] | 1 | 1.2 | - |

### A. Experiment Settings

*1) Datasets:* We used five representative datasets in crowd scenes, i.e., WiderFace [12] for dense face detection, ShanghaiTech [21] for crowd counting and localization, NWPU-Crowd [22] for crowd localization, and CARPK and PUCPR+ [38] for vehicle counting from the drone view.

**WiderFace** is one of the most challenging face detection benchmarks, where the 32,203 images contain 393,703 human-labeled faces were captured in a wide variety of imaging conditions, such as large variations in scale and pose, high occlusion, and changeable illumination conditions. 40%, 10%, and 50% of the images were used for training, validation, and testing, respectively. Following existing point-supervised detection methods [9], [10], we trained our model on the training set, and reported detection and counting results on the validation set.

**ShanghaiTech** presents high-density crowds, which contains 482 images on Part_A (SHA) and 716 images on Part_B (SHB). The number of people in an image ranges from 33 to 3139 on SHA, and 9 to 578 on SHB. We followed the training and testing split in [21] to evaluate the counting and central point localization performance.

**NWPU-Crowd** provides a large-scale benchmark for crowd counting and localization. It consists of 3109, 500, and 1500 images for the training, validation, and test, respectively, and the images contain more than 2 million annotated heads with points and boxes.

**CARPK and PUCPR+** are composed of images of parking lots from the drone view and high-rise buildings, respectively. CARPK contains nearly 90k cars, while PUCPR+ contains about 17k cars in total. We used the evaluation protocol in their benchmark to evaluate the counting performance of our method.

*2) Implementation details:* The backbone, ResNet-50, was initialized with the pre-trained weights on ImageNet [45]. We trained our detector on 3 GPUs with the batch size of 12. We adopted Adam [46] optimizer with the learning rate of $7.5 \times 10^{-6}$ for the WiderFace dataset and $7.5 \times 10^{-5}$ for the remaining ones. The input images were randomly re-scaled, color distorted, flipped, and then cropped into $704 \times 704$ image patches. We used the same re-scale technique as CSP [20]. All training processes can only access to the point annotations. For the datasets with bounding box annotations, i.e., WiderFace, CARPK, and PUCPR+, we calculated their center points for training. We stopped training at 200k, 8k, 60k, 45k, and 4.5k iterations for WiderFace, ShanghaiTech, NWPU-

TABLE II
COMPARISONS OF OUR LUDA-BASED PSEUDO SIZE GENERATION METHOD WITH THE GAK-BASED METHOD ON THE TRAINING SET OF WIDERFACE [12]. REFINEMENT OF PSEUDO OBJECT SIZES BY OUR SELF-TRAINING APPROACH IS LISTED AT THE LAST ENTRY. THE NUMBER IN THE PARENTHESES DENOTES THE IoU THRESHOLD FOR AP CALCULATION. THE LARGER THE VALUE OF AP IS, THE BETTER.

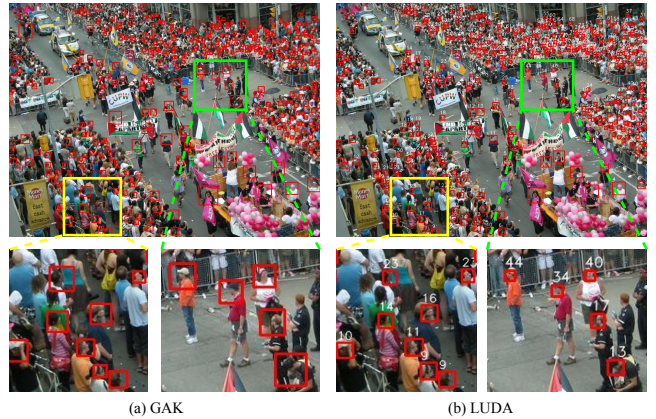| Size Generation | AP (0.3) | AP (0.5) | AP (0.7) |
|---|---|---|---|
| GAK [21] | 55.8 | 29.5 | 7.0 |
| LUDA (ours) | 60.3 | 31.2 | 7.5 |
| Refined size (ours) | **79.4** | **50.6** | **14.4** |



(a) GAK　　　　　　　　　(b) LUDA

Fig. 3. Visual comparison of the pseudo object sizes generated by GAK and LUDA in a crowded scene. For LUDA, the number on top of each bounding box refers to the object crowdedness. The high crowdedness value means the pseudo size is accurate.

Crowd, CARPK, and PUCPR+, respectively. Like [9], [20], we performed the multi-scale testing to generate bounding boxes. Then, non-maximum suppression (NMS) was used to filter the generated boxes. Following [10], we randomly took 10% training images of ShanghaiTech as the validation set. For ShanghaiTech and NWPU-Crowd, we chose the best model by performing a threshold search to minimize the counting error on the validation set. The other datasets adopted the confidence threshold of 0.4 to produce boxes for counting.

As mentioned in Sec. III-B, we experimentally chose multiple sets of parameters by which the precise pseudo sizes were generated for crowded objects on training sets. Table I shows the parameter settings used in Eqs. (3) and (7) for different datasets. We especially set the maximum value of $\alpha_j$ to 50 to avoid large gradients since dense crowds of ShanghaiTech produce large values of $\alpha_j$. For all datasets, we set $\eta = 1$ in Eq. (7).

*3) Evaluation metrics:* For the detection task, we adopted the evaluation protocol in WiderFace [12] to calculate average precision (AP). The true positive is defined as the intersection of union (IoU) between ground truth boxes and detected boxes greater than a threshold of $0.5$. For the counting task, we adopted the commonly-used mean absolute error (MAE) and root mean square error (RMSE) to evaluate the distance between the predicted counts and the ground-truth ones. The MAE indicates the accuracy of methods, while the RMSE reflects their robustness. They are defined as

$$\text{MAE} = \frac{1}{N_t} \sum_{N_t}^{1} |\hat{c}_i - c_i|, \text{RMSE} = \sqrt{\frac{1}{N_t} \sum_{N_t}^{1} (\hat{c}_i - c_i)^2},$$
(8)

where $N_t$ is the total number of testing images, and $\hat{c}_i$ and $c_i$ are the estimated count and ground-truth count of the $i$-th image, respectively. When counting on WiderFace, we used a normalized MAE (NAE) [14], which normalizes the absolute error by the ground-truth face count.

For the center point localization task, we adopted two evaluation metrics, i.e., AP and mean localization error (MLE) respectively from [9] and [10]. AP defines true positive center point as those whose distance to its ground truth is smaller than a threshold of 20 pixels. MLE calculates the distances in pixels between the predicted center points and the ground truth, and then averages the distances over the testing set. One-to-one matching associates the predictions and the ground truth. The lower the value of MLE is, the better. The AP and MLE are

suitable and reasonable for the datasets without bounding-box annotations, e.g., ShanghaiTech. For NWPU-Crowd, we used its evaluation protocol [22], i.e., the Precision, Recall, and F1-measure. These metrics evaluate the localization ability of detection-based algorithms in crowded environments.

### B. Ablation Study

We conducted ablation studies on the WiderFace [12] benchmark to evaluate and analyze the improvement of several important modules of our approach, including the locally-uniform distribution assumption (LUDA), the crowdedness-aware loss, and the confidence and order-aware refinement scheme. We evaluated our approach on both the detection and counting tasks. Table II shows the AP results of pseudo size generation methods on the training set of WiderFace, while Table III shows the AP, MAE, and NAE results on the validation set of WiderFace. The results of state-of-the-art point-supervised detection method (i.e., PSDDN [9]) and the crowed counting method (i.e., Shi *et al.* [14]) are respectively shown in the 1st and 2nd entries of Table III for comparison.

*1) Effectiveness of the pseudo object size generation method:* We compared the proposed LUDA-based pseudo size generation method with the geometry-adaptive kernel (GAK) based method in [9], [21]. The generated boxes were evaluated on the training set of WiderFace. We calculated the AP in three IoU thresholds of 0.3, 0.5, and 0.7, so the true positives become gradually harder to reach. The results are shown in Table II, where it can be observed that the proposed LUDA-based method obtains higher AP scores under all thresholds, compared with the GAK-based method, validating the advantage of our LUDA. Such an advantage also means that we can generate more accurate pseudo bounding boxes at the beginning of training.

Fig. 3 shows a visual comparison of the pseudo object sizes generated by GAK and LUDA in a crowded scene. We can observe that: 1) both methods generate more accurate sizes for crowded objects than sparse objects; and 2) LUDA makes the sparse objects' sizes be accurate when they are surrounded by

TABLE III

THE ABLATIVE STUDIES TOWARDS THE CROWDEDNESS-AWARE LOSS AND THE CONFIDENCE AND ORDER-AWARE REFINEMENT SCHEME, AS WELL AS COMPARISONS WITH STATE-OF-THE-ART POINT-SUPERVISED DETECTION METHOD [9] AND COUNTING METHOD [14] ON THE VALIDATION SET OF WIDERFACE [12]. "C-BY-D" MEANS COUNTING BY DETECTION, AND "C-BY-R" MEANS COUNTING BY REGRESSION. "SIZE INIT." STANDS FOR DIRECTLY APPLYING THE PSEUDO OBJECT SIZE INITIALIZATION METHOD ON THE VALIDATION SET TO OBTAIN THE DETECTION RESULT. "-" IS NOT APPLICABLE. THE LARGER THE VALUE OF AP IS, THE BETTER. THE LOWER THE VALUES OF MAE AND NAE ARE, THE BETTER.

| Method | Category | Crowdedness-aware loss | Confidence and order-aware refinement scheme | Detection AP | | | Counting | |
|---|---|---|---|---|---|---|---|---|
| | | | | easy | medium | hard | MAE | NAE |
| PSDDN [9] | C-by-D | - | - | 60.5 | 60.5 | 39.6 | - | - |
| Shi *et al.* [14] | C-by-R | - | - | - | - | - | 3.2 | 0.40 |
| GAK [21] | Size Init. | - | - | 7.1 | 12.5 | 27.3 | - | - |
| LUDA (ours) | Size Init. | - | - | 7.2 | 12.8 | 29.5 | - | - |
| Ours | C-by-D | ✗ | ✗ | 18.6 | 23.4 | 30.8 | 3.4 | 0.81 |
| | | ✓ | ✗ | 21.3 | 26.8 | 35.2 | 3.6 | 0.97 |
| | | ✗ | ✓ | 55.1 | 55.1 | 52.5 | 2.3 | **0.27** |
| | | ✓ | ✓ | **75.8** | **71.0** | **64.4** | **2.2** | 0.29 |

crowded objects, as illustrated in the green and yellow boxes of Fig. 3.

We also evaluated the GAK-based and LUDA-based methods by yielding bounding boxes on the validation set of WiderFace. The AP scores are shown in the 3rd and 4th entries of Table III, where it can be seen that the accuracy of the generated bounding boxes is extremely low, and especially the AP score is only 7.2% on the easy subset. This observation shows that only applying pseudo size generation methods cannot obtain accurate object sizes even with the ground-truth center points.

*2) Effectiveness of the crowdedness-aware loss and the confidence and order-aware refinement scheme:* The crowdedness-aware loss and the confidence and order-aware refinement scheme are the critical components to improve the AP for point-supervised detection. Based on the pseudo bounding boxes, we trained the detector on the training set of WiderFace with four sittings 1) using only the crowdedness-aware loss, 2) using only the confidence and order-aware refinement scheme, 3) using neither of the two modules, and 4) using both of the two modules. The AP, MAE, and NAE results are listed in the 5th to 8th entries of Table III.

**Pseudo sizes only.** The detector was trained with the pseudo object sizes only, as listed in the 5th entry of Table III. The AP scores are improved by about 11% on the easy and medium subsets in comparison with the LUDA-based size generation method. However, the values are still low, e.g., 18.6%, 23.4%, and 30.8% AP on the easy, medium, and hard subsets, respectively. This observation demonstrates that it is hard to attain an acceptable detector when trained only with pseudo bounding boxes.

**Either the crowdedness-aware loss or the confidence and order-aware refinement scheme.** As shown in the 6th and 7th entries of Table III, with the use of the crowdedness-aware loss (resp. the confidence and order-aware refinement scheme), the AP scores increase to 21.3% (easy), 26.8% (medium), and 30.8% (hard) (resp. 55.1% (easy), 55.1% (medium), and 52.5% (hard)), which validate the effectiveness of these two modules. Moreover, it can be known that the refinement module is more effective than the loss. Note that only with the refinement scheme, our approach outperforms PSDDN [9] on the hard subset in the detection task (52.5% AP vs. 39.4% AP) and Shi *et al.* [14] in the counting task (2.3 MAE vs. 3.2 MAE).
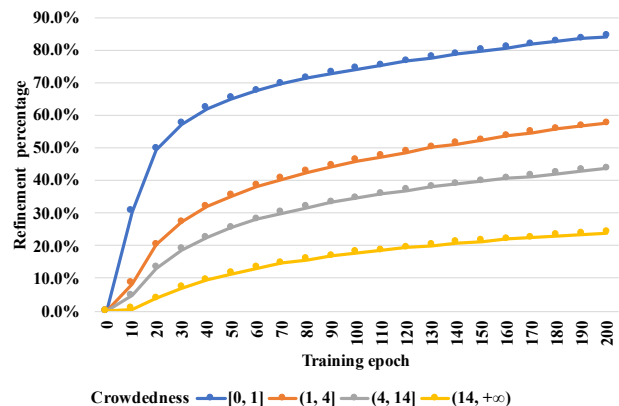


Fig. 4. The percentage of refined training examples with respect to the training epoch.

**Both the crowdedness-aware loss and the confidence and order-aware refinement scheme.** When both modules are activated, our approach achieves the best results in both detection and counting tasks. There is a big jump of AP from using a single module to both modules. We think that the crowdedness-aware loss emphasizes the accurate pseudo sizes, but neglects to update the noisy sizes during training, while the confidence and order-aware refinement scheme updates the pseudo sizes but it overlooks the importance of the accurate pseudo sizes in the network's weights updating. The combination of such two modules could well compensate to each other. Numerically, the AP scores of our approach with both two modules increase to 75.8% (easy), 71.0% (medium), and 64.4% (hard). Compared with PSDDN [9], our method improves the AP by more than 10%. The MAE and NAE of our approach respectively decrease approximately 31.2% and 27.5% against [14]. The results demonstrate the significant superiority of our method in both detection and counting tasks.

Besides, the last row of Table II shows that the refined pseudo sizes improve the AP (0.3) and AP (0.5) by nearly 20%. The gradually enhanced quality of training examples helps the detector become stronger. To demonstrate this improvement, we plot some training examples before and after training on the WiderFace dataset in Fig. 1. Especially for large and sparse objects, the bounding boxes are refined to
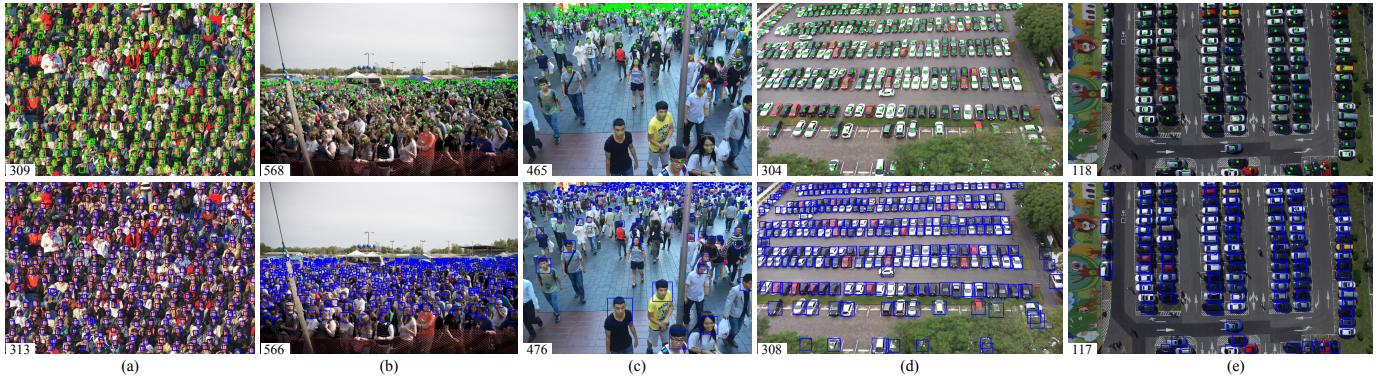
Fig. 5. Qualitative results on (a) WiderFace [12], (b) SHA [21], (c) SHB [21], (d) PUCPR+ [38], and (e) CARPK [38]. The top row shows the ground-truth boxes or points, and counts. The bottom row shows the bounding boxes and counts predicted by our approach. *Zoom in the figure for better viewing.*

<div style="display:flex">
<div>

TABLE IV
COMPARISONS WITH STATE-OF-THE-ART POINT-SUPERVISED DETECTION METHODS ON THE VALIDATION SET OF WIDERFACE [12]. "*": THE AP RESULTS ARE PROVIDED BY [9]. THE LARGER THE VALUE OF AP IS, THE BETTER. THE INITIALIZED PSEUDO BOXES WERE USED WHEN TRAINING CSP WITH POINT SUPERVISION.

| Method | Supervision | easy | medium | hard |
|---|---|---|---|---|
| LSC-CNN [10] | Box | 57.3 | 70.1 | 68.9 |
| Faceness-WIDER [47] | Box | 71.3 | 63.4 | 34.5 |
| Faster R-CNN* [11] | Box | 84.0 | 72.4 | 34.7 |
| CSP (anchor-free) [20] | Box | 90.7 | **95.2** | **96.1** |
| HR-TinyFace [48] | Box | **92.5** | 91.0 | 80.6 |
| CSP (anchor-free) [20] | Point | 18.6 | 23.4 | 30.8 |
| CSR-A-thr [10] | Point | 30.2 | 41.9 | 33.5 |
| PSDNN [9] | Point | 60.5 | 60.5 | 39.6 |
| LSC-CNN [10] | Point | 40.5 | 62.1 | 46.2 |
| Ours | Point | **75.8** | **71.0** | **64.4** |

</div>
<div>

TABLE V
COMPARISONS WITH CROWD LOCALIZATION METHODS ON SHA AND SHB [21] DATASETS. THE LARGER THE VALUE OF AP IS AND THE LOWER THE VALUE OF MLE IS, THE BETTER.

| Method | SHA AP (%) | SHA MLE | SHB AP (%) | SHB MLE |
|---|---|---|---|---|
| PSDDN [9] | 73.7 | - | 75.9 | - |
| CSR-A-thr [10] | - | 16.8 | - | 12.3 |
| LSC-CNN [10] | 67.6 | 9.6 | 76.0 | 9.0 |
| Ours | **85.3** | **8.0** | **91.6** | **6.0** |

ods, including PSDDN [9] and LSC-CNN [10], and the counting-by-regression method, i.e., CSR-A-thr [10]. The CSR-A-thr is the detection version of CSRNet [13]. Table IV shows the AP scores of different methods on the validation set of WiderFace [12]. From Table IV, we can see that our method outperforms the other point-supervised methods to a significant margin (i.e., more than 10% AP improvement). Especially on the hard subset, our method improves the AP score of the second best LSC-CNN method by 18.2%.

Besides, Table IV provides the results of several box-supervised methods for comparison, where it can be seen that our method even outperforms box-supervised LSC-CNN and Faceness-WIDER on the easy subset and Faceness-WIDER and Faster-RCNN on the hard subset. Although the AP score of the box-supervised CSP achieves above 90%, it significantly drops to 30.8% when trained with point-initialized pseudo boxes, which demonstrates the effectiveness of our point-supervised training scheme. Fig. 5 shows some visual results of bounding box predicted by our method.

*2) Center point localization of crowds:* **SHA and SHB datasets.** To evaluate the localization ability for the datasets with only point-level annotations, we compared our method with PSDNN [9], CSR-A-thr [10], and LSC-CNN [10] on the SHA and SHB [21] datasets. Table V presents the results of the AP and MLE metrics. Our approach obtains the best results, showing 85.3% AP and 8.0 MLE for SHA, and 91.6% AP and 6.0 MLE for SHB.

**NWPU-Crowd dataset.** We also compared our approach with state-of-the-art crowd localization methods on the large-scale NWPU-Crowd [22] dataset. The results are shown in Table VI, where it can be observed that our method obtains

</div>
</div>

encompass the face regions.

*3) Behavior of the confidence and order-aware refinement scheme:* In our confidence and order-aware refinement scheme, the pseudo object sizes are updated in a robust and orderly manner, i.e., the easy examples (sparse and large objects) with inaccurate sizes are first refined, followed by the hard examples (crowded and small objects). The object sizes are refined when their posterior probabilities are larger than the threshold of 0.6. To illustrate the behavior of such a refinement scheme, we counted the percentage of examples whose probabilities exceed 0.6 during training. Specifically, we split examples into four categories corresponding to the crowdedness intervals $[0, 1]$, $(1, 4]$, $(4, 14]$, and $(14, +\infty)$ with 64k, 30k, 31k, and 30k training examples, respectively. Fig. 4 shows the percentage of examples with their posterior probabilities larger than 0.6 under each category. We can observe that the percentage of the revised low-crowdedness examples increases more rapidly than that of the high-crowdedness examples in the first 20 epochs. Throughout the entire training epochs, the percentage corresponding to lower-crowdedness reaches the highest, demonstrating the effectiveness of our refinement scheme.

### C. Comparison with State-of-the-Art Methods

*1) Point-supervised face detection:* We compared our approach with state-of-the-art point-supervised detection meth-
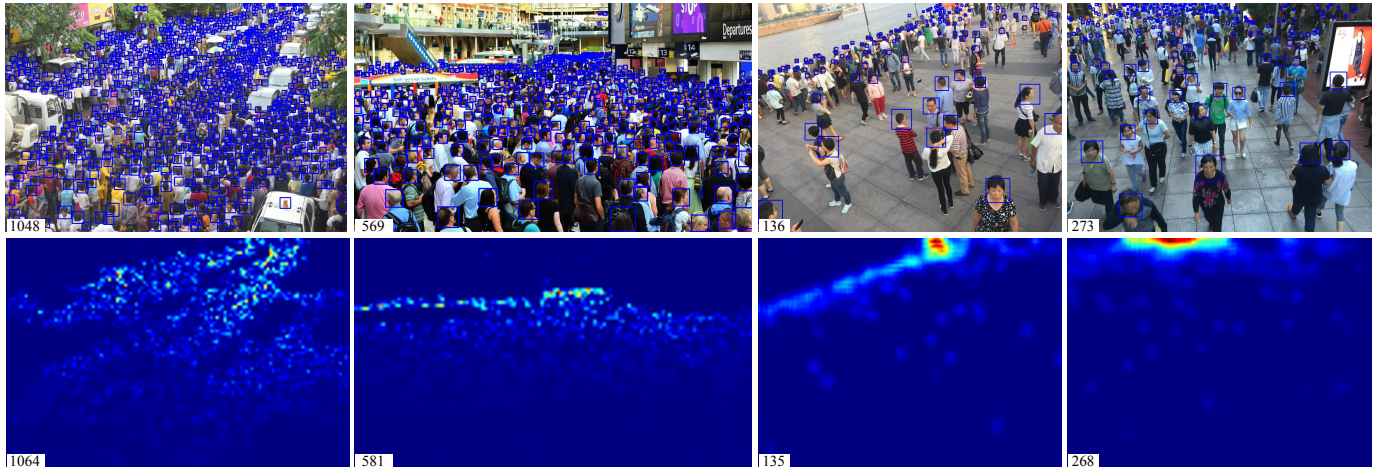
Fig. 6. Comparisons between the bounding boxes (at top row) produced by our approach and the density maps (at bottom row) produced by CSRNet [13] on the SHA and SHB [21] datasets. The predicted count is shown at the bottom-left corner of an image. The ground-truth counts from left to right images are 1068, 584, 139, and 274. *Zoom in the figure for better viewing.*

TABLE VI
COMPARISON WITH STATE-OF-THE-ART CROWD LOCALIZATION METHODS ON THE NWPU-CROWD TEST SET. THE VALUES OF F1-MEASURE, PRECISION, AND RECALL WERE CALCULATED UNDER THE THRESHOLD $\sigma_l$ [22].

| Method | Supervision | Output | F1 (%) | Precision (%) | Recall (%) |
|--------|-------------|--------|--------|---------------|------------|
| Faster R-CNN [11] | Box | Box | 6.7 | **95.8** | 3.5 |
| TinyFaces [48] | Box | Box | 56.7 | 52.9 | 61.1 |
| VGG+GPR [49] | Point | Point | 52.5 | 55.8 | 49.6 |
| RAZ_Loc [50] | Point | Point | 59.8 | 66.6 | 54.3 |
| Ours | Point | Box | **63.7** | 65.1 | **62.4** |

TABLE VII
COMPARISONS WITH STATE-OF-THE-ART COUNTING-BY-REGRESSION METHODS (IN THE TOP PART) AND COUNTING-BY-DETECTION METHODS (IN THE BOTTOM PART) ON SHA AND SHB [21], AND WIDERFACE [12] (WF) DATASETS. "-": THE AUTHOR DOES NOT PROVIDE THE RESULT. "*": THE RESULTS ARE PROVIDED BY [10]. THE LOWER THE VALUES OF MAE AND RMSE ARE, THE BETTER.

| Method | SHA MAE | SHA RMSE | SHB MAE | SHB RMSE | WF MAE |
|--------|---------|----------|---------|----------|--------|
| Zhang *et al.* [21] | 110.2 | 173.2 | 26.4 | 41.3 | 7.1 |
| CSRNet [13] | 68.2 | 115.0 | 10.6 | 16.0 | 4.3 |
| Cao et al. [33] | 67.0 | 104.5 | 8.4 | 13.6 | 8.5 |
| PSDNN+ [9] | 65.9 | 112.3 | 9.1 | 14.2 | - |
| Shi *et al.* [14] | 65.2 | 109.4 | **7.2** | **12.2** | **3.2** |
| PSD+DCL [35] | 65.0 | 108.0 | 8.1 | 13.3 | - |
| HA-CCN [37] | **62.9** | **94.9** | 8.1 | 13.4 | **-** |
| TinyFace* [48] | 237.8 | 422.8 | - | - | - |
| LC-FCN8 [28] | - | - | 13.1 | - | - |
| PSDNN [9] | 85.4 | 159.2 | 16.1 | 27.9 | - |
| LSC-CNN [10] | 66.4 | 117.0 | 8.1 | 15.7 | - |
| Ours | **65.1** | **104.4** | **7.8** | **12.6** | **2.2** |

TABLE VIII
COMPARISONS WITH STATE-OF-THE-ART VEHICLE COUNTING METHODS ON THE CARPK AND PUCPR+ BENCHMARK [38]. THE LOWER THE VALUES OF MAE AND RMSE ARE, THE BETTER.

| Method | CARPK MAE | CARPK RMSE | PUCPR+ MAE | PUCPR+ RMSE |
|--------|-----------|------------|------------|-------------|
| LPN Counting [38] | 23.80 | 36.79 | 22.76 | 34.46 |
| RetinaNet [23] | 16.62 | 22.30 | 24.58 | 33.12 |
| IEP Counting [39] | 51.83 | - | 15.17 | - |
| Goldman *et al.* [26] | 6.77 | 8.52 | 7.16 | 12.00 |
| Li *et al.* [40] | 5.24 | 7.38 | 3.92 | 5.06 |
| Ours | **4.95** | **7.09** | **3.20** | **4.83** |

*3) Crowd counting:* In addition to detection, we evaluated our method in the crowd counting task. Table VII shows the results on SHA and SHB [21], and WiderFace [12] datasets. For a fair comparison, we split the counting methods into two categories: counting-by-regression methods and counting-by-detection methods. Our method achieves the best performance when compared with state-of-the-art counting-by-detection methods. For SHA and SHB, the proposed method is comparable to state-of-the-art counting-by-regression methods, such as Shi *et al.* [14] and HA-CCN [37].

Note that our method not only provides the count but also estimates the bounding boxes for object instances. We qualitatively evaluated the proposed method by visualizing the predicted bounding boxes on the SHA, and SHB datasets in Fig. 5. Besides, Fig. 6 shows the predicted boxes by our method and the estimated density maps by CSRNet [13] on the SHA and SHB datasets. We argue that the box outputs of our method are more informative than the density maps of counting-by-regression methods because the boxes provide high-level understanding of crowds.

*4) Vehicle counting from drone view:* To evaluate the generalization ability in other domains, we trained our model on the CARPK and PUCPR+ datasets [38]. The datasets provide densely-packed car counting images from drone view. Table VIII reports the MAE and RMSE values of our method and state-of-the-art vehicle counting methods [23], [26], [38]–

the best F1-measure (i.e., 63.7%) and Recall (i.e., 62.4%). Fast R-CNN achieves 95.8% Precision but sacrifices the Recall to 3.5%, indicating that Fast R-CNN fails to detect crowded people. Note that the proposed approach is the only one that can output bounding boxes when trained with point supervision.

[40]. It can be seen that our method consistently performs better than the other methods, demonstrating that our model is flexible for various detection and counting tasks. Some visual results on CARPK and PUCPR+ are shown in Fig. 5.
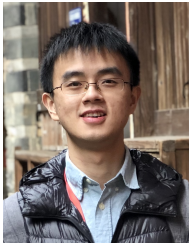
## V. CONCLUSION

In this paper, we have presented a self-training approach to train a typical detector with only point-level annotations such that the detector can accurately detect and count objects in crowd scenes simultaneously. This is achieved by the locally-uniform distribution assumption, the crowdedness-aware loss, the confidence and order-aware refinement scheme, and the decoding method to promote the detector to generate accurate bounding boxes in a coarse-to-fine and end-to-end manner. Extensive experimental results over multiple commonly used benchmark datasets have demonstrated that the proposed approach achieves the best performance in point-supervised detection and counting tasks among detection-based methods, and our method even produces comparable performance to state-of-the-art counting-by-regression methods. We believe that DNN-based object detection in crowds with only point supervision is a potential and promising research area.

## REFERENCES

[1] V. A. Sindagi and V. M. Patel, "A survey of recent advances in cnn-based single image crowd counting and density estimation," *Pattern Recognit. Lett.*, vol. 107, pp. 3–16, 2018.

[2] T. Li, H. Chang, M. Wang, B. Ni, R. Hong, and S. Yan, "Crowded scene analysis: A survey," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 25, no. 3, pp. 367–386, 2015.

[3] D. Kang, Z. Ma, and A. B. Chan, "Beyond counting: comparisons of density maps for crowd analysis tasks-counting, detection, and tracking," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 29, no. 5, pp. 1408–1422, 2019.

[4] A. B. Chan, Z.-S. J. Liang, and N. Vasconcelos, "Privacy preserving crowd monitoring: Counting people without people models or tracking," in *CVPR*. IEEE, 2008, pp. 1–7.

[5] Y.-L. Chen, B.-F. Wu, H.-Y. Huang, and C.-J. Fan, "A real-time vision system for nighttime vehicle detection and traffic surveillance," *IEEE Trans. Ind. Electron.*, vol. 58, no. 5, pp. 2030–2044, 2011.

[6] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *ECCV*. Springer, 2014, pp. 740–755.

[7] X. Wang, J. Chen, Z. Wang, W. Liu, S. Satoh, C. Liang, and C.-W. Lin, "When pedestrian detection meets nighttime surveillance: A new benchmark," in *IJCAI*, 2020, pp. 509–515.

[8] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: A review," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 11, pp. 3212–3232, 2019.

[9] Y. Liu, M. Shi, Q. Zhao, and X. Wang, "Point in, box out: Beyond counting persons in crowds," in *CVPR*, 2019, pp. 6469–6478.

[10] D. B. Sam, S. V. Peri, M. Narayanan Sundararaman, A. Kamath, and R. V. Babu, "Locate, size and count: Accurately resolving people in dense crowds via detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.

[11] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015, pp. 91–99.

[12] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "Wider face: A face detection benchmark," in *CVPR*, 2016, pp. 5525–5533.

[13] Y. Li, X. Zhang, and D. Chen, "Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes," in *CVPR*. IEEE, 2018, pp. 1091–1100.

[14] Z. Shi, P. Mettes, and C. G. Snoek, "Counting with focus for free," in *ICCV*, 2019, pp. 4200–4209.

[15] G. Gao, J. Gao, Q. Liu, Q. Wang, and Y. Wang, "Cnn-based density estimation and crowd counting: A survey," *arXiv preprint arXiv:2003.12783*, 2020.

[16] M. Rodriguez, I. Laptev, J. Sivic, and J.-Y. Audibert, "Density-aware person detection and tracking in crowds," in *ICCV*. IEEE, 2011, pp. 2423–2430.

[17] Y. Xiang, A. Alahi, and S. Savarese, "Learning to track: Online multi-object tracking by decision making," in *ICCV*, 2015, pp. 4705–4713.

[18] M. A. Turk and A. P. Pentland, "Face recognition using eigenfaces," in *CVPR*. IEEE, 1991, pp. 586–591.

[19] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *CVPR*, 2014, pp. 152–159.

[20] W. Liu, S. Liao, W. Ren, W. Hu, and Y. Yu, "High-level semantic feature detection: A new perspective for pedestrian detection," in *CVPR*, June 2019, pp. 5187–5196.

[21] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma, "Single-image crowd counting via multi-column convolutional neural network," in *CVPR*. IEEE, 2016, pp. 589–597.

[22] Q. Wang, J. Gao, W. Lin, and X. Li, "Nwpu-crowd: A large-scale benchmark for crowd counting," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.

[23] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollar, "Focal loss for dense object detection," in *ICCV*, Oct 2017, pp. 2980–2988.

[24] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," *arXiv preprint arXiv:1904.07850*, 2019.

[25] K. Duan, S. Bai, L. Xie, H. Qi, Q. Huang, and Q. Tian, "Centernet: Keypoint triplets for object detection," in *CVPR*, 2019, pp. 6569–6578.

[26] E. Goldman, R. Herzig, A. Eisenschtat, J. Goldberger, and T. Hassner, "Precise detection in densely packed scenes," in *CVPR*, 2019, pp. 5227–5236.

[27] Q. Cai, Y. Pan, Y. Wang, J. Liu, T. Yao, and T. Mei, "Learning a unified sample weighting network for object detection," in *CVPR*, 2020, pp. 14 173–14 182.

[28] I. H. Laradji, N. Rostamzadeh, P. O. Pinheiro, D. Vazquez, and M. Schmidt, "Where are the blobs: Counting by localization with point supervision," in *ECCV*, 2018, pp. 547–562.

[29] H. Idrees, M. Tayyab, K. Athrey, D. Zhang, S. Al-Maadeed, N. Rajpoot, and M. Shah, "Composition loss for counting, density map estimation and localization in dense crowds," in *ECCV*, 2018, pp. 532–546.

[30] C. Zhang, H. Li, X. Wang, and X. Yang, "Cross-scene crowd counting via deep convolutional neural networks," in *CVPR*. IEEE, 2015, pp. 833–841.

[31] Y. Wang, J. Hou, and L.-P. Chau, "Object counting in video surveillance using multi-scale density map regression," in *ICASSP*. IEEE, 2019, pp. 2422–2426.

[32] A. B. Chan and N. Vasconcelos, "Bayesian poisson regression for crowd counting," in *ICCV*. IEEE, 2009, pp. 545–551.

[33] X. Cao, Z. Wang, Y. Zhao, and F. Su, "Scale aggregation network for accurate and efficient crowd counting," in *ECCV*, 2018, pp. 734–750.

[34] V. Lempitsky and A. Zisserman, "Learning to count objects in images," in *NeurIPS*, 2010, pp. 1324–1332.

[35] Q. Wang, W. Lin, J. Gao, and X. Li, "Density-aware curriculum learning for crowd counting," *IEEE Trans. Cybern.*, pp. 1–13, 2020.

[36] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.

[37] V. A. Sindagi and V. M. Patel, "Ha-ccn: Hierarchical attention-based crowd counting network," *IEEE Trans. Image Process.*, vol. 29, pp. 323–335, 2019.

[38] M.-R. Hsieh, Y.-L. Lin, and W. H. Hsu, "Drone-based object counting by spatially regularized regional proposal network," in *ICCV*, 2017, pp. 4145–4153.

[39] T. Stahl, S. L. Pintea, and J. C. van Gemert, "Divide and count: Generic object counting by image divisions," *IEEE Trans. Image Process.*, vol. 28, no. 2, pp. 1035–1044, 2018.

[40] W. Li, H. Li, Q. Wu, X. Chen, and K. N. Ngan, "Simultaneously detecting and counting dense vehicles from drone images," *IEEE Trans. Ind. Electron.*, vol. 66, no. 12, pp. 9651–9662, 2019.

[41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.

[42] R. Girshick, "Fast r-cnn," in *ICCV*. IEEE, Dec 2015, pp. 1440–1448.

[43] A. Shrivastava, A. Gupta, and R. Girshick, "Training region-based object detectors with online hard example mining," in *CVPR*, 2016, pp. 761–769.

[44] A. RoyChowdhury, P. Chakrabarty, A. Singh, S. Jin, H. Jiang, L. Cao, and E. Learned-Miller, "Automatic adaptation of object detectors to new domains using self-training," in *CVPR*, 2019, pp. 780–790.

[45] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NeurIPS*, 2012, pp. 1097–1105.

[46] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[47] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "From facial parts responses to face detection: A deep learning approach," in *ICCV*, 2015, pp. 3676–3684.

[48] P. Hu and D. Ramanan, "Finding tiny faces," in *CVPR*, 2017, pp. 951–959.

[49] J. Gao, T. Han, Q. Wang, and Y. Yuan, "Domain-adaptive crowd counting via inter-domain features segregation and gaussian-prior reconstruction," *arXiv preprint arXiv:1912.03677*, 2019.

[50] C. Liu, X. Weng, and Y. Mu, "Recurrent attentive zooming for joint crowd counting and precise localization," in *CVPR*, 2019, pp. 1217–1226.

**Yi Wang** received the B.Eng. degree in electronic information engineering and M.Eng. degree in information and signal processing from the School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China, in 2013 and 2016, respectively. He is currently a research associate with the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, and he is also working toward the Ph.D. degree of Nanyang Technological University. His research interests include image restoration, image recognition, object detection, and crowd analysis.
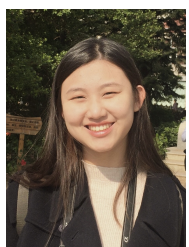
**Junhui Hou** received the B.Eng. degree in information engineering (Talented Students Program) from the South China University of Technology, Guangzhou, China, in 2009, the M.Eng. degree in signal and information processing from Northwestern Polytechnical University, Xi'an, China, in 2012, and the Ph.D. degree in electrical and electronic engineering from the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, in 2016. He has been an Assistant Professor with the Department of Computer Science, City University of Hong Kong, since 2017. His research interests fall into the general areas of visual computing, such as image/video/3D geometry data representation, processing and analysis, semi/un-supervised data modeling, and data compression and adaptive transmission

Dr. Hou was the recipient of several prestigious awards, including the Chinese Government Award for Outstanding Students Study Abroad from China Scholarship Council in 2015, and the Early Career Award (3/381) from the Hong Kong Research Grants Council in 2018. He is a member of MSA-TC and VSPC-TC, IEEE CAS. He is currently serving as an Associate Editor for IEEE Transactions on Circuits and Systems for Video Technology, The Visual Computer, and Signal Processing: Image Communication, and the Guest Editor for the IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing. He also served as an Area Chair of ACM MM 2019 and 2020, IEEE ICME 2020, and WACV 2021. He is a senior member of IEEE.

**Xinyu Hou** received the B.Eng. degree from the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, in 2020. She is currently a project officer at the School of Computer Science and Engineering, Nanyang Technological University. Her research interests include image/video generation, computer vision, and machine learning.

**Lap-Pui Chau** received the Bachelor degree from Oxford Brookes University, and the Ph.D. degree from The Hong Kong Polytechnic University, in 1992 and 1997, respectively. He is Assistant Chair (Academic) of School of Electrical and Electronic Engineering, Nanyang Technological University. His research interests include visual signal processing algorithms, light-field imaging, video analytics for intelligent transportation system, and human motion analysis.

He was a General Chairs for IEEE International Conference on Digital Signal Processing (DSP 2015) and International Conference on Information, Communications and Signal Processing (ICICS 2015). He was a Program Chairs for Visual Communications and Image Processing (VCIP 2020 and VCIP 2013), International Conference on Digital Signal Processing (DSP 2018), International Conference on Multimedia and Expo (ICME 2016) and International Symposium on Intelligent Signal Processing and Communications Systems (ISPACS 2010).

He was the chair of Technical Committee on Circuits & Systems for Communications (TC-CASC) of IEEE Circuits and Systems Society from 2010 to 2012. He served as an associate editor for IEEE Transactions on Multimedia, IEEE Signal Processing Letters, IEEE Transactions on Circuits and Systems for Video Technology, IEEE Transactions on Circuits and Systems II, and is currently serving as an associate editor for IEEE Transactions on Broadcasting, and The Visual Computer (Springer Journal). Besides, he was an IEEE Distinguished Lecturer for 2009-2019, and a steering committee member of IEEE Transactions for Mobile Computing from 2011-2013. He is an IEEE Fellow.