# Object Counting and Instance Segmentation with Image-level Supervision

Hisham Cholakkal *        Guolei Sun *        Fahad Shahbaz Khan        Ling Shao

Inception Institute of Artificial Intelligence, UAE

`firstname.lastname@inceptioniai.org`

## Abstract

*Common object counting in a natural scene is a challenging problem in computer vision with numerous real-world applications. Existing image-level supervised common object counting approaches only predict the global object count and rely on additional instance-level supervision to also determine object locations. We propose an image-level supervised approach that provides both the global object count and the spatial distribution of object instances by constructing an object category density map. Motivated by psychological studies, we further reduce image-level supervision using a limited object count information (up to four). To the best of our knowledge, we are the first to propose image-level supervised density map estimation for common object counting and demonstrate its effectiveness in image-level supervised instance segmentation. Comprehensive experiments are performed on the PASCAL VOC and COCO datasets. Our approach outperforms existing methods, including those using instance-level supervision, on both datasets for common object counting. Moreover, our approach improves state-of-the-art image-level supervised instance segmentation [33] with a relative gain of 17.8% in terms of average best overlap, on the PASCAL VOC 2012 dataset.* [1]

## 1. Introduction

Common object counting, also referred as *generic object counting*, is the task of accurately predicting the number of different object category instances present in natural scenes (see Fig. 1). The common object categories in natural scenes can vary from fruits to animals and the counting must be performed in both indoor and outdoor scenes (*e.g.* COCO or PASCAL VOC datasets). Existing works employ a localization-based strategy or utilize regression-based models directly optimized to predict object count,
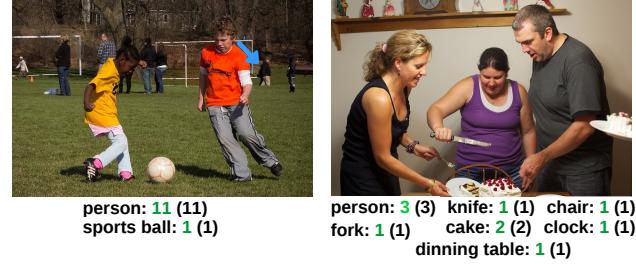
---

Figure 1. Object counting on COCO dataset. The ground-truth and our predictions are shown in black and green, respectively. Despite being trained using image-level object counts within the subitizing range [1-4], it accurately counts objects beyond the subitizing range (11 persons) under heavy occlusion (marked with blue arrow to show two persons) in the left image and diverse object categories in the right.

where the latter has been shown to provide superior results [4]. However, regression-based methods only predict the global object count without determining object locations. Beside global counts, the spatial distribution of objects in the form of a per-category density map is helpful in other tasks, e.g., to delineate adjacent objects in instance segmentation (see Fig. 2).

The problem of density map estimation to preserve the spatial distribution of people is well studied in crowd counting [3, 16, 18, 22, 31]. Here, the global count for the image is obtained by summing over the predicted density map. Standard crowd density map estimation methods are required to predict large number of person counts in the presence of occlusions, e.g., in surveillance applications. The key challenges of constructing a density map in natural scenes are different to those in crowd density estimation, and include large intra-class variations in generic objects, co-existence of multiple instances of different objects in a scene (see Fig. 1), and sparsity due to many objects having zero count on multiple images.

Most methods for crowd density estimation use instance-level (point-level or bounding box) supervision that requires manual annotation of each instance location. Image-

(a) Input Image      (b) PRM [33]      (c) Our Approach      (d) Our Density Map
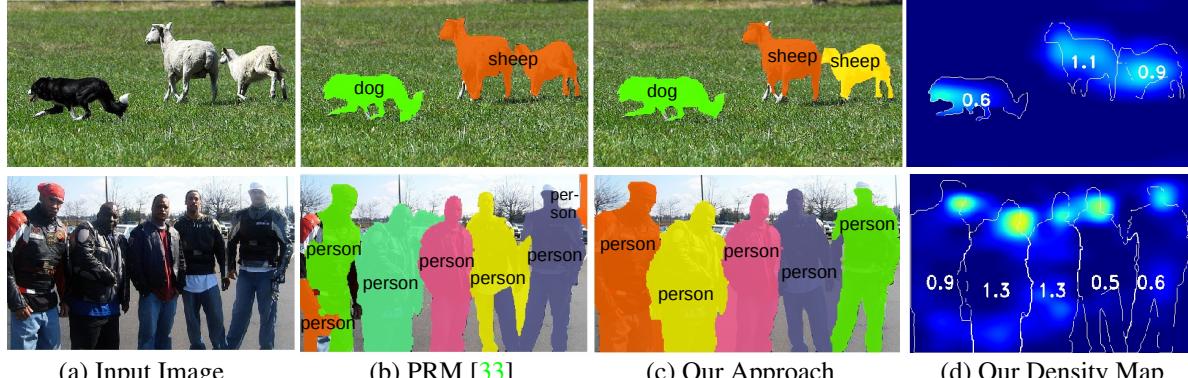
Figure 2. Instance segmentation examples using the PRM method [33] (b) and our approach (c), on the PASCAL VOC 2012. Top row: The PRM approach [33] fails to delineate spatially adjacent two sheep category instances. Bottom row: single person parts predicted as multiple persons along with inaccurate mask separation results in over-prediction (7 instead of 5). Our approach produces accurate masks by exploiting the spatial distribution of object count in per-category density maps (d). Density map accumulation for each predicted mask is shown inside the contour drawn for clarity. In the top row, density maps for sheep and dog categories are overlaid.

level supervised training alleviates the need for such user-intensive annotation by requiring only the count of different object instances in an image. We propose an image-level supervised density map estimation approach for natural scenes, that predicts the global object count while preserving the spatial distribution of objects.

Even though image-level supervised object counting reduces the burden of human annotation and is much weaker compared to instance-level supervisions, it still requires each object instance to be counted sequentially. Psychological studies [2, 6, 12, 20] have suggested that humans are capable of counting objects non-sequentially using holistic cues for fewer object counts, termed as a subitizing range (generally 1-4). We utilize this property to further reduce image-level supervision by only using object count annotations within the subitizing range. For short, we call this image-level lower-count (ILC) supervision. Chattopadhyay *et al*. [4] also investigate common object counting, where object counts (both within and beyond the subitizing range) are used to predict the global object count. Alternatively, instance-level (bounding box) supervision is used to count objects by dividing an image into non-overlapping regions, assuming each region count falls within the subitizing range. Different to these strategies [4], our ILC supervised approach requires neither bounding box annotation nor information beyond the subitizing range to predict both the count and the spatial distribution of object instances.

In addition to common object counting, the proposed ILC supervised density map estimation is suitable for other scene understanding tasks. Here, we investigate its effectiveness for image-level supervised instance segmentation, where the task is to localize each object instance with pixel-level accuracy, provided image-level category labels. Recent work of [33], referred as peak response map (PRM), tackles the problem by boosting the local maxima (peaks)

in the class response maps [23] of an image classifier using a peak stimulation module. A scoring metric is then used to rank off-the-shelf object proposals [21, 25] corresponding to each peak for instance mask prediction. However, PRM struggles to delineate spatially adjacent object instances from the same object category (see Fig. 2(b)). We introduce a penalty term into the scoring metric that assigns a higher score to object proposals with a predicted count of 1, providing improved results (Fig. 2(c)). The predicted count is obtained by accumulating the density map over the entire object proposal region (Fig. 2(d)).

**Contributions:** We propose an ILC supervised density map estimation approach for common object counting. A novel loss function is introduced to construct per-category density maps with explicit terms for predicting the global count and spatial distribution of objects. We also demonstrate the applicability of the proposed approach for image-level supervised instance segmentation. For common object counting, our ILC supervised approach outperforms state-of-the-art instance-level supervised methods with a relative gain of 6.4% and 2.9%, respectively, in terms of mean root mean square error (mRMSE), on the PASCAL VOC 2007 and COCO datasets. For image-level supervised instance segmentation, our approach improves the state of the art from 37.6 to 44.3 in terms of average best overlap (ABO), on the PASCAL VOC 2012 dataset.

## 2. Related work

Chattopadhyay *et al*. [4] investigated regression-based common object counting, using image-level (per-category count) and instance-level (bounding box) supervisions. The image-level supervised strategy, denoted as glancing, used count annotations from both within and beyond the subitizing range to predict the global count of objects, without providing information about their location. The instance-level
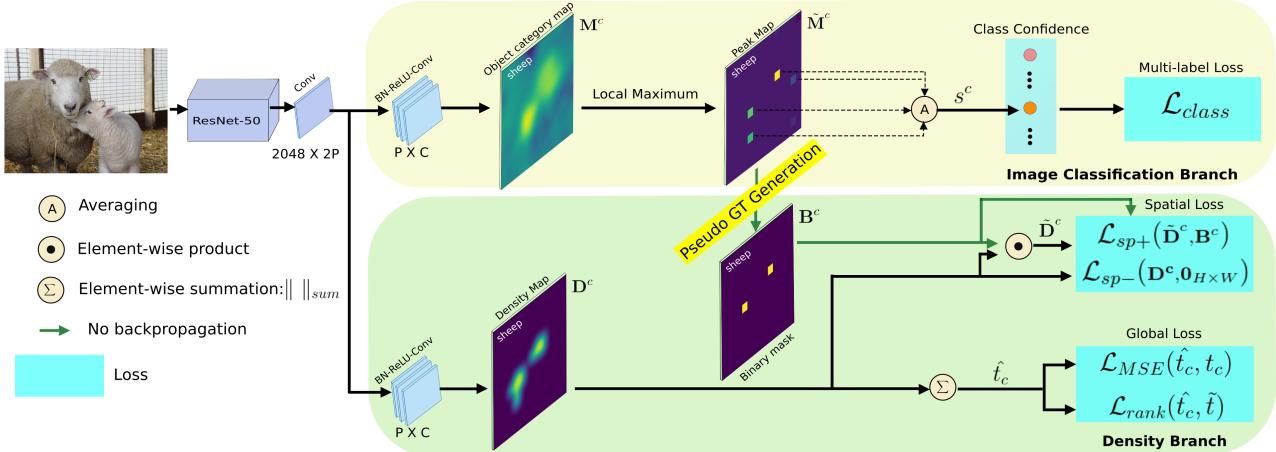
Figure 3. Overview of our overall architecture. Our network has an image classification and a density branch, trained jointly using ILC supervision. The image classification branch predicts the presence and absence of objects. This branch is used to generate pseudo ground-truth for training the density branch. The density branch has two terms (spatial and global) in the loss function and produces a density map to predict the global object count and preserve the spatial distribution of objects.

(bounding box) supervised strategy, denoted as subitizing, estimated a large number of objects by dividing an image into non-overlapping regions, assuming the object count in each region falls within the subitizing range. Instead, our ILC supervised approach requires neither bounding box annotation nor beyond subitizing range count information during training. It then predicts the global object count, even beyond the subitizing range, together with the spatial distribution of object instances.

Recently, Laradji *et al.* [14] proposed a localization-based counting approach, trained using instance-level (point) supervision [1]. During inference, the model outputs blobs indicating the predicted locations of objects of interest and uses [29] to estimate object counts from these blobs. Different to [14], our approach is image-level supervised and directly predicts the object count through a simple summation of the density map without requiring any post-processing [29]. Regression-based methods generally perform well in the presence of occlusions [4, 15], while localization-based counting approaches [9, 14] generalize well with a limited number of training images [14, 15]. Our method aims to combine the advantages of both approaches through a novel loss function that jointly optimizes the network to predict object locations and global object counts in a density map.

Reducing object count supervision for salient object subitizing was investigated in [30]. However, their task is class-agnostic and subitizing is used to only count within the subitizing range. Instead, our approach constructs category-specific density maps and accurately predicts object counts both within and beyond the subitizing range. Common object counting has been previously used to improve object detection [4, 8]. Their approach only uses the count information during detector training with no explicit

component for count prediction. In contrast, our approach explicitly learns to predict the global object count.

## 3. Proposed method

Here, we present our image-level lower-count (ILC) supervised density map estimation approach. Our approach is built upon an ImageNet pre-trained network backbone (ResNet50) [11]. The proposed network architecture has two output branches: image classification and density branch (see Fig. 3). The image classification branch estimates the presence or absence of objects, whereas the density branch predicts the global object count and the spatial distribution of object instances by constructing a density map. We remove the global pooling layer from the backbone and adapt the fully connected layer with a $1 \times 1$ convolution having $2P$ channels as output. We divide these $2P$ channels equally between the image classification and density branches. We then add a $1 \times 1$ convolution having $C$ output channels in each branch, resulting in a fully convolutional network [19]. Here, $C$ is the number of object categories and $P$ is empirically set to be proportional to $C$. In each branch, the convolution is preceded by a batch normalization and a ReLU layer. The first branch provides object category maps and the second branch produces a density map for each object category.

### 3.1. The Proposed Loss Function

Let $\mathbf{I}$ be a training image and $\mathbf{t} = \{t_1, t_2, ..., t_c, ..., t_C\}$ be the corresponding vector for the ground-truth count of $C$ object categories. Instead of using an absolute object count, we employ a lower-count strategy to reduce the amount of image-level supervision. Given an image $\mathbf{I}$, object categories are divided into three non-overlapping sets based on their respective instance counts. The first set,

$A$, indicates object categories which are absent in $\mathbf{I}$ (i.e., $t_c = 0$). The second set, $S$, represents categories within the subitizing range (i.e, $0 < t_c \leq 4$). The final set, $\tilde{S}$, indicates categories beyond the subitizing range (i.e, $t_c \geq \tilde{t}$, where $\tilde{t} = 5$).

Let $\mathbf{M} = \{\mathbf{M}^1, \mathbf{M}^2, ..., \mathbf{M}^c, ..., \mathbf{M}^C\}$ denote the object category maps in the image classification branch, where $\mathbf{M}^c \in R^{H \times W}$. Let $\mathbf{D} = \{\mathbf{D}^1, \mathbf{D}^2, ..., \mathbf{D}^c, ..., \mathbf{D}^C\}$ represent density maps produced by the density branch, where $\mathbf{D}^c \in R^{H \times W}$. Here, $H \times W$ is the spatial size of both the object category and density maps. The image classification and density branches are jointly trained, in an end-to-end fashion, given only ILC supervision with the following loss function:

$$\mathcal{L} = \mathcal{L}_{class} + \underbrace{\mathcal{L}_{spatial} + \mathcal{L}_{global}}_{Density\ map\ branch}. \qquad (1)$$

Here, the first term refers to multi-label image classification loss [13] (see Sec. 3.1.1). The last two terms, $\mathcal{L}_{spatial}$ and $\mathcal{L}_{global}$, are used to train the density branch (Sec. 3.1.2).

### 3.1.1  Image Classification Branch

Generally, training a density map requires instance-level supervision, such as point-level annotations [15]. Such information is unavailable in our ILC supervised setting. To address this issue, we propose to generate pseudo ground-truth by exploiting the coarse-level localization capabilities of an image classifier [23, 32] via object category maps. These object category maps are generated from a fully convolutional architecture shown in Fig. 3.

While specifying classification confidence at each image location, class activation maps (CAMs) struggle to delineate multiple instances from the same object category [23, 32]. Recently, the local maxima of CAMs are further boosted, to produce object category maps, during an image-classifier training for image-level supervised instance segmentation [33]. Boosted local maxima aim at falling on distinct object instances. For details on boosting local maxima, we refer to [33]. Here, we use local maxima locations to generate pseudo ground-truth for training the density branch.

As described earlier, object categories in $\mathbf{I}$ are divided into three non-overlapping sets: $A$, $S$ and $\tilde{S}$. To train a one-versus-rest image classifier, we derive binary labels from $t_c$ that indicate the presence $\forall c \in \{S, \tilde{S}\}$ or absence $\forall c \in A$ of object categories. Let $\tilde{\mathbf{M}}^c \in R^{H \times W}$ be the peak map derived from $c^{th}$ object category map ($\mathbf{M}^c$) of $\mathbf{M}$ such that:

$$\tilde{\mathbf{M}}^c(i,j) = \begin{cases} \mathbf{M}^c(i,j), & \text{if } \mathbf{M}^c(i,j) > \mathbf{M}^c(i-r_i, j-r_j), \\ 0, & \text{otherwise.} \end{cases}$$

Here, $-r \leq r_i \leq r$, $-r \leq r_j \leq r$ where $r$ is the radius for the local maxima computation. We set $r = 1$, as in [33]. The local maxima are searched at all spatial locations with a stride of one. To train an image classifier, a class confidence

score $s^c$ of the $c^{th}$ object category is computed as the average of non-zero elements of $\tilde{\mathbf{M}}^c$. In this work, we use the multi-label soft-margin loss [13] for binary classification.

### 3.1.2  Density Branch

The image classification branch described above predicts the presence or absence of objects by using the class confidence scores derived from the peak map $\tilde{\mathbf{M}}^c$. However, it struggles to differentiate between multiple objects and single object parts due to the lack of prior information about the number of object instances (see Fig. 2(b)). This causes a large number of false positives in the peak map $\tilde{\mathbf{M}}^c$. Here, we utilize the count information and introduce a pseudo ground-truth generation scheme that prevents training a density map at those false positive locations.

When constructing a density map, it is desired to estimate accurate object counts at any image sub-region. Our spatial loss term $\mathcal{L}_{spatial}$ in Eq. 1 ensures that individual object instances are localized while the global term $\mathcal{L}_{global}$ constrains the global object count to that of the ground-truth. This enables preservation of the spatial distribution of object counts in a density map. Later, we show that this property helps to improve instance segmentation.

**Spatial Loss:** The spatial loss $\mathcal{L}_{spatial}$ is divided into the loss $\mathcal{L}_{sp+}$ which enhances the positive peaks corresponding to instances of object categories within $S$, and the loss $\mathcal{L}_{sp-}$ which suppresses false positives of categories within $A$. Due to the unavailability of absolute object count, the set $\tilde{S}$ is not used in the spatial loss and treated separately later. To enable ILC supervised density map training using $\mathcal{L}_{spatial}$, we generate a pseudo ground-truth binary mask from peak map $\tilde{\mathbf{M}}^c$.

*Pseudo Ground-truth Generation:* To compute the spatial loss $\mathcal{L}_{sp+}$, a pseudo ground-truth is generated for set $S$. For all object categories $c \in S$, the $t_c$-th highest peak value of peak map $\tilde{M}^c$ is computed using the heap-max algorithm [5]. The $t_c$-th highest peak value $h_c$ is then used to generate a pseudo ground-truth binary mask $\mathbf{B}^c$ as,

$$\mathbf{B}^c = u(\tilde{\mathbf{M}}^c - h_c). \qquad (2)$$

Here, $u(n)$ is the unit step function which is 1 only if $n \geq 0$. Although the non-zero elements of the pseudo ground-truth mask $\mathbf{B}^c$ indicate object locations, its zero elements do not necessarily point towards the background. Therefore, we construct a masked density map $\tilde{\mathbf{D}}^c$ to exclude density map $\mathbf{D}^c$ values at locations where the corresponding $\mathbf{B}^c$ values are zero. Those density map $\mathbf{D}^c$ values should also be excluded during the loss computation in Eq. 4 and backpropagation (see Sec. 3.2), due to their risk of introducing false negatives. This is achieved by computing the Hadamard product between the density map $\mathbf{D}^c$ and $\mathbf{B}^c$ as,

$$\tilde{\mathbf{D}}^c = \mathbf{D}^c \odot \mathbf{B}^c. \qquad (3)$$

The spatial loss $\mathcal{L}_{sp+}$ for object categories within the subitizing range $S$ is computed between $\mathbf{B}^c$ and $\tilde{\mathbf{D}}^c$ using a logistic binary cross entropy (logistic BCE) [24] loss for positive ground-truth labels. The logistic BCE loss transfers the network prediction ($\tilde{\mathbf{D}}^c$) through a sigmoid activation layer $\sigma$ and computes the standard BCE loss as,

$$\mathcal{L}_{sp+}(\tilde{\mathbf{D}}^c, \mathbf{B}^c) = -\sum_{\forall c \in S} \frac{\|\mathbf{B}^c \odot \log(\sigma(\tilde{\mathbf{D}}^c))\|_{sum}}{|S| \cdot \|\mathbf{B}^c\|_{sum}}. \quad (4)$$

Here, $|S|$ is the cardinality of the set $S$ and the norm $\| \; \|_{sum}$ is computed by taking the summation over all elements in a matrix. For example, $\|\mathbf{B}^c\|_{sum} = \mathbf{1}^h \mathbf{B}^c \mathbf{1}^w$, where $\mathbf{1}^h$ and $\mathbf{1}^w$ are all-ones vectors of size $1 \times H$ and $W \times 1$, respectively. Here, the highest $t_c$ peaks in $\tilde{\mathbf{M}}^c$ are assumed to fall on $t_c$ instances of object category $c \in S$. Due to the unavailability of ground-truth object locations, we use this assumption and observe that it holds in most scenarios.

The spatial loss $\mathcal{L}_{sp+}$ for the positive ground-truth labels enhances positive peaks corresponding to instances of object categories within $S$. However, the false positives of the density map for $c \in S$ are not penalized in this loss. We therefore introduce another term, $\mathcal{L}_{sp-}$, into the loss function to address the false positives of $c \in A$. For $c \in A$, positive activations of $\mathbf{D}^c$ indicate false detections. A zero-valued mask $\mathbf{0}_{H \times W}$ is used as ground-truth to reduce such false detections using logistic BCE loss,

$$\mathcal{L}_{sp-}(\mathbf{D}^c, \mathbf{0}_{H \times W}) = -\sum_{c \in A} \frac{\|\log(1 - \sigma(\mathbf{D}^c)\|_{sum}}{|A| \cdot H \cdot W}. \quad (5)$$

Though the spatial loss ensures the preservation of spatial distribution of objects, only relying on local information may result in deviations in the global object count.

**Global Loss:** The global loss penalizes the deviation of the predicted count $\hat{t}_c$ from the ground-truth. It has two components: ranking loss $\mathcal{L}_{rank}$ for object categories beyond the subitizing range (i.e., $\forall c \in \tilde{S}$) and mean-squared error (MSE) loss $\mathcal{L}_{MSE}$ for the rest of the categories. $\mathcal{L}_{MSE}$ penalizes the predicted density map, if the global count prediction does not match with the ground-truth count. i.e.,

$$\mathcal{L}_{MSE}(\hat{t}_c, t_c) = \sum_{c \in \{A, S\}} \frac{(\hat{t}_c - t_c)^2}{|A| + |S|}. \quad (6)$$

Here, the predicted count $\hat{t}_c$ is the accumulation of the density map for a category $c$ over its entire spatial region. i.e. $\hat{t}_c = \|\mathbf{D}^c\|_{sum}$. Note that object categories in $\tilde{S}$ were not previously considered in the computation of spatial loss $\mathcal{L}_{spatial}$ and mean-squared error loss $\mathcal{L}_{MSE}$. Here, we introduce a ranking loss [28] with a zero margin that penalizes under-counting for object categories within $\tilde{S}$,

$$\mathcal{L}_{rank}(\hat{t}_c, \tilde{t}) = \sum_{c \in \tilde{S}} \frac{max(0, \tilde{t} - \hat{t}_c)}{|\tilde{S}|}. \quad (7)$$

The ranking loss penalizes the density branch if the predicted object count $\hat{t}_c$ is less than $\tilde{t}$ for $c \in \tilde{S}$. Recall, the beyond subitizing range $\tilde{S}$ starts from $\tilde{t} = 5$.

Within the subitizing range $S$, the spatial loss term $\mathcal{L}_{spatial}$ is optimized to locate object instances while the global MSE loss ($\mathcal{L}_{MSE}$) is optimized for accurately predicting the corresponding global count. Due to the joint optimization of both these terms within the subitizing range, the network learns to correlate between the located objects and the global count. Further, the network is able to locate object instances, generalizing beyond the subitizing range $\tilde{S}$ (see Fig. 2). Additionally, the ranking loss $\mathcal{L}_{rank}$ term in the proposed loss function ensures the penalization of under counting beyond the subitizing range $\tilde{S}$.

**Mini-batch Loss:** Normalized loss terms $\hat{\mathcal{L}}_{sp+}$, $\hat{\mathcal{L}}_{sp-}$, $\hat{\mathcal{L}}_{MSE}$ and $\hat{\mathcal{L}}_{rank}$ are computed by averaging respective loss terms over all images in the mini-batch. The $\mathcal{L}_{spatial}$ is computed by $\hat{\mathcal{L}}_{sp+} + \hat{\mathcal{L}}_{sp-}$. For categories beyond the subitizing range, $\hat{\mathcal{L}}_{rank}$ can lead to over-estimation of the count. Hence, $\mathcal{L}_{global}$ is computed by assigning a relatively lower weight ($\lambda = 0.1$) to $\hat{\mathcal{L}}_{rank}$ (see Table. 2). i.e., $\mathcal{L}_{global} = \hat{\mathcal{L}}_{MSE} + \lambda * \hat{\mathcal{L}}_{rank}$.

### 3.2. Training and Inference

Our network is trained in two stages. In the first stage, the density branch is trained with only $\mathcal{L}_{MSE}$ and $\mathcal{L}_{rank}$ losses using $S$ and $\tilde{S}$ respectively. The spatial loss $\mathcal{L}_{spatial}$ in Eq. 1 is excluded in the first stage, since it requires a pseudo ground-truth generated from the image classification branch. The second stage includes the spatial loss.

**Backpropagation:** We use $\mathbf{B}^c$ derived from the image classification branch as a pseudo ground-truth to train the density branch. Therefore, the backproapation of gradients through $\mathbf{B}^c$ to the classifier branch is not required (shown with green arrows in Fig. 3). The image classification branch is backpropagated as in [33]. In the density branch, we use Hadamard product of the density map with $\mathbf{B}^c$ in Eq. 3 to compute $\mathcal{L}_{sp+}$ for $c \in S$. Hence, the gradients ($\delta^c$) for the $c^{th}$ channel of the last convolution layer of the density branch, due to $\mathcal{L}_{sp+}$, are computed as,

$$\delta_{sp+}^c = \frac{\partial \hat{\mathcal{L}}_{sp+}}{\partial \tilde{\mathbf{D}}^c} \odot \mathbf{B}^c. \quad (8)$$

Since $\mathcal{L}_{MSE}$, $\mathcal{L}_{rank}$ and $\mathcal{L}_{sp-}$ are computed using MSE, ranking and logistic BCE losses on convolution outputs, their respective gradients are computed using off-the-shelf pytorch implementation [24].

**Inference:** The image classification branch outputs a class confidence score $s^c$ for each class, indicating the presence ($\hat{t}_c > 0$, if $s^c > 0$) or absence ($\hat{t}_c = 0$, if $s^c \leq 0$) of the object category $c$. The predicted count $\hat{t}_c$ is obtained by summing the density map $\mathbf{D}^c$ for category $c$ over its entire spatial region. The proposed approach only utilizes subitizing

annotations ($t_c \leq 4$) and accurately predicts object counts for *both* within and beyond subitizing range (see Fig. 6).

### 3.3. Image-level Supervised Instance Segmentation

The proposed ILC supervised density map estimation approach can also be utilized for instance segmentation. Note that the local summation of an ideal density map over a ground-truth segmentation mask is one. We use this property to improve state-of-the-art image-level supervised instance segmentation (PRM) [33]. PRM employs a scoring metric that combines instance level cues from peak response maps $R$, class aware information from object category maps and spatial continuity priors from off-the-shelf object proposals [21, 25]. Here, the peak response maps are generated from local maxima (peaks of $\tilde{\mathbf{M}}^c$) through a peak back-propagation process [33]. The scoring metric is then used to rank object proposals corresponding to each peak for instance mask prediction. We improve the scoring metric by introducing an additional term $d_p$ in the metric. The term $d_p$ penalizes an object proposal $P_r$, if the predicted count in those regions of the density map $\mathbf{D}^c$ is different from one, as $d_p = |1 - \|\mathbf{D}^c \cdot P_r\|_{sum}|$. Here, $|\ |$ is the absolute value operator. For each peak, the new scoring metric $Score$ selects the highest scoring object proposal $P_r$.

$$Score = \alpha \cdot R * P_r + R * \hat{P}_r - \beta \cdot Q * P_r - \gamma \cdot d_p. \quad (9)$$

Here, the background mask $Q$ is derived from object category map and $\hat{P}_r$ is the contour mask of the proposal $P_r$ derived using morphological gradient [33]. Parameters $\alpha$, $\beta$ and $\gamma$ are empirically set as in [33].

## 4. Experiments

**Implementation details:** Throughout our experiments, we fix the training parameters. An initial learning rate of $10^{-4}$ is used for the pre-trained ResNet-50 backbone, while image classification and density branches are trained with an initial learning rate of 0.01. The number of input channels $P$ of $1 \times 1$ convolution for each branch is set to $P = 1.5 \times C$. A mini-batch size of 16 is used for the SGD optimizer. The momentum is set to 0.9 and weight decay to $10^{-4}$. Considering high imbalance between non-zero and zero counts in COCO dataset (*e.g.* 79 negative categories for each positive category), only 10% of samples in the set $A$ are used to train the density branch. Code will be made public upon publication.

**Datasets:** We evaluate common object counting on the PASCAL VOC 2007 [7] and COCO [17] datasets. For fair comparison, we employ same splits, named as count-train, count-val and count-test, as used in the state-of-the-art methods [14], [4]. For COCO dataset, the training set is used as count-train, first half of the validation set as the count-val and its second half as the count-test. In Pascal VOC 2007 dataset, we evaluated against the count of non-difficult instances in the count-test as in [14]. For instance

| Approach | SV | mRMSE | mRMSE-nz | m-relRMSE | m-relRMSE-nz |
|---|---|---|---|---|---|
| CAM+regression | IC | 0.45 | 1.52 | 0.29 | 0.64 |
| Peak+regression | IC | 0.64 | 2.51 | 0.30 | 1.06 |
| Proposed | ILC | **0.29** | **1.14** | **0.17** | **0.61** |

Table 1. Counting performance on the Pascal VOC 2007 count-test set using our approach and two baselines. Both baselines are obtained by training the network using the MSE loss function.



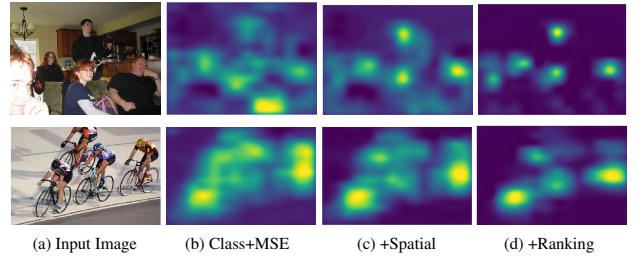(a) Input Image    (b) Class+MSE    (c) +Spatial    (d) +Ranking

Figure 4. Progressive improvement in density map quality with the incremental introduction of spatial and ranking loss terms. In both cases (top row: person and bottom row: bicycle), our overall loss function integrating all three terms provides the best density maps. The global object count is accurately predicted (top row: 5 persons and bottom row: 4 bicycles) by accumulation of the respective density map.

segmentation, we train and report the results on the PASCAL VOC 2012 dataset similar to [33].

**Evaluation Criteria:** The predicted count $\hat{t}_c$ is rounded to the nearest integer. We evaluate common object counting, as in [4, 14], using root squared error (RMSE) metric and its three variants namely RMSE non-zero (RMSE-nz), relative RMSE (relRMSE) and relative RMSE non-zero (relRMSE-nz). The $RMSE_c$ and $relRMSE_c$ errors for category $c$ are computed as $\sqrt{\frac{1}{T} \sum_{i=1}^{T} (t_{ic} - \hat{t_{ic}})^2}$ and, $\sqrt{\frac{1}{T} \sum_{i=1}^{T} \frac{(t_{ic} - \hat{t_{ic}})^2}{t_{ic}+1}}$ respectively. Here, $T$ is the total number of images in the test set and $\hat{t}_{ic}$, $t_{ic}$ are the predicted and ground-truth counts for image $i$. The errors are then averaged across all categories to obtain the mRMSE and m-relRMSE on a dataset. The above metrics are also evaluated for ground-truth instances with non-zero counts as mRMSE-nz and m-relRMSE-nz. For all error metrics, smaller numbers indicate better performance. We refer to [4] for more details. For instance segmentation, the performance is evaluated using Average Best Overlap (ABO) [26] and $mAP^r$, as in [33]. The $mAP^r$ is computed with intersection over union (IoU) thresholds of 0.25, 0.5 and 0.75.

**Supervision Levels:** The level of supervision is indicated as SV in Tab. 3 and 4. BB indicates bounding box supervision and PL indicates point-level supervision for each object instance. Image-level supervised methods using only within subitizing range counts are denoted as ILC, while the methods using both within and beyond subitizing range counts are indicated as IC.

| | $\mathcal{L}_{class}+$ $\mathcal{L}_{MSE}$ | $\mathcal{L}_{class}+$ $\mathcal{L}_{spatial}$ $+\mathcal{L}_{MSE}$ | $\mathcal{L}$ $\lambda=0.1$ | $\mathcal{L}$ $\lambda=0.01$ | $\mathcal{L}$ $\lambda=0.05$ | $\mathcal{L}$ $\lambda=0.5$ | $\mathcal{L}$ $\lambda=1$ |
|---|---|---|---|---|---|---|---|
| mRMSE | 0.36 | 0.33 | **0.29** | 0.31 | 0.30 | 0.32 | 0.36 |
| mRMSE-nz | 1.52 | 1.32 | **1.14** | 1.27 | 1.16 | 1.23 | 1.40 |

Table 2. Left: Progressive integration of different terms in loss function and its impact on the final counting performance on the PASCAL VOC count-test set. Right: influence of the weight ($\lambda$) of ranking loss.

## 4.1. Common Object Counting Results

**Ablation Study:** We perform an ablation study on the PASCAL VOC 2007 count-test. First, the impact of our two-branch architecture is analyzed by comparing it with two baselines: class-activation [32] based regression (CAM+regression) and peak-based regression (Peak+regression) using the local-maximum boosting approach of [33]. Both baselines are obtained by end-to-end training of the network, employing the same backbone, using MSE loss function to directly predict global count. Tab. 1 shows the comparison. Our approach largely outperforms both baseline highlighting the importance of having a two-branch architecture with explicit terms in the loss function to preserve the spatial distribution of objects. Next, we evaluate the contribution of each term in our loss function towards the final count performance.

Fig. 4 shows the systematic improvement in density maps (top row: person and bottom row: bicycle) quality with the incremental addition of (c) spatial $\mathcal{L}_{spatial}$ and (d) ranking ($\mathcal{L}_{rank}$) loss terms to the (b) MSE ($\mathcal{L}_{rank}$) loss term. Similar to CAM, the density branch trained with MSE loss alone gives coarse location of object instances. However, many background pixels are identified as part of the object (false positives) resulting in inaccurate spatial distribution of object instances. Further, this inclusion of false positives prevents the delineation of multiple object instances. Incorporating the spatial loss term improves the spatial distribution of objects in both density maps. The density maps are further improved by the incorporation of the ranking term that penalizes the under-estimation of count beyond the subitizing range (top row) in the loss function. Moreover, it also helps to reduce the false positives within the subitizing range (bottom row). Tab. 2 shows the systematic improvement, in terms of mRMSE and mRMSE-nz, when integrating different terms in our loss function. The best results are obtained when integrating all three terms (classification, spatial and global) in our loss function. We also evaluate the influence of $\lambda$ that controls the relative weight of the ranking loss. We observe $\lambda = 0.1$ provides the best results and fix it for all datasets.

**State-of-the-art Comparison:** Tab. 3 and 4 show state-of-the-art comparisons for common object counting on the PASCAL VOC 2007 and COCO datasets respectively. On the PASCAL VOC 2007 dataset (Tab. 3), the glancing approach (glance-noft-2L) of [4] using image-level

| Approach | SV | mRMSE | mRMSE-nz | m-relRMSE | m-relRMSE-nz |
|---|---|---|---|---|---|
| Aso-sub-ft-3×3 [4] | BB | 0.43 | 1.65 | 0.22 | 0.68 |
| Seq-sub-ft-3×3 [4] | BB | 0.42 | 1.65 | 0.21 | 0.68 |
| ens [4] | BB | 0.42 | 1.68 | 0.20 | 0.65 |
| Fast-RCNN [4] | BB | 0.50 | 1.92 | 0.26 | 0.85 |
| LC-ResFCN [14] | PL | 0.31 | 1.20 | 0.17 | 0.61 |
| LC-PSPNet [14] | PL | 0.35 | 1.32 | 0.20 | 0.70 |
| glance-noft-2L [4] | IC | 0.50 | 1.83 | 0.27 | 0.73 |
| Proposed | ILC | **0.29** | **1.14** | **0.17** | **0.61** |

Table 3. State-of-the-art counting performance comparison on the Pascal VOC 2007 count-test. Our ILC supervised approach outperforms existing methods.

| Approach | SV | mRMSE | mRMSE-nz | m-relRMSE | m-relRMSE-nz |
|---|---|---|---|---|---|
| Aso-sub-ft-3×3 [4] | BB | 0.38 | 2.08 | 0.24 | 0.87 |
| Seq-sub-ft-3×3 [4] | BB | 0.35 | 1.96 | 0.18 | 0.82 |
| ens [4] | BB | 0.36 | 1.98 | 0.18 | **0.81** |
| Fast-RCNN [4] | BB | 0.49 | 2.78 | 0.20 | 1.13 |
| LC-ResFCN [14] | PL | 0.38 | 2.20 | 0.19 | 0.99 |
| glance-ft-2L [4] | IC | 0.42 | 2.25 | 0.23 | 0.91 |
| Proposed | ILC | **0.34** | **1.89** | **0.18** | 0.84 |

Table 4. State-of-the-art counting performance comparison on the COCO count-test set. Despite using reduced supervision, our approach provides superior results compared to existing methods on three metrics. Compared to the image-level count (IC) supervised approach [4], our method achieves an absolute gain of 8% in terms of mRMSE.



orange: 2, 8 (8)  carrot: 2, 5 (5)  bowl: 0, 1 (1)  person: 4, 1 (1)  broccoli: 1, 5 (5)  zebra: 15, 12 (12)  person: 5, 6 (6)  remote: 2, 1 (1)  tv: 1, 1 (1)

Figure 5. Object counting examples on the COCO dataset. The ground-truth, point-level supervised counts [14] and our predictions are shown in black, red and green respectively. Our approach accurately performs counting beyond the subitizing range and on diverse categories (fruits to animals) under heavy occlusions (highlighted by a red arrow in the left image).

supervision both within and beyond the subitizing range (IC) achieves mRMSE score of $0.50$. Our ILC supervised approach considerably outperforms the glance-noft-2L method with a absolute gain of 21% in mRMSE. Furthermore, our approach achieves consistent improvements on all error metrics, compared to state-of-the-art point-level and bounding box based supervised methods.

Tab. 4 shows the results on COCO dataset. Among the existing methods, the two BB supervised approaches (Seq-sub-ft-3x3 and ens) yields mRMSE scores of $0.35$ and $0.36$ respectively. The PL supervised LC-ResFCN approach [14] achieves mRMSE score of $0.38$. The IC supervised glancing approach (glance-noft-2L) obtains mRMSE score of $0.42$. Our approach outperforms the glancing approach with an absolute gain of 8% in mRMSE. Furthermore, our approach also provides consistent improvements over the glancing approach in the other three error metrics and is only below the two BB supervised methods (Seq-sub-ft3x3 and ens) in m-relRMSE-nz. Fig. 5 shows object counting examples using our approach and the point-level
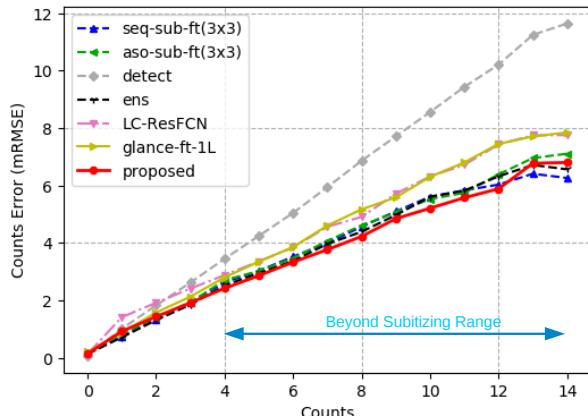
Figure 6. Counting performance comparison in RMSE, across all categories, at different ground-truth count values on the COCO count-test set. Different methods, including BB and PL supervision, are shown in the legend. Our ILC supervised approach provides superior results compared to the image-level supervised glancing method. Furthermore, our approach performs favourably compared to other methods using instance-level supervision.



(a) Input Image    (b) PRM [33]    (c) Our Approach

Figure 7. Instance segmentation examples obtained using PRM [33] and our approach. The proposed approach accurately delineates spatially adjacent multiple object instances of horse and cow categories.

| Method | $mAP_{0.25}^r$ | $mAP_{0.5}^r$ | $mAP_{0.75}^r$ | ABO |
|---|---|---|---|---|
| MELM+MCG [27] | 36.9 | 22.9 | 8.4 | 32.9 |
| CAM+MCG [32] | 20.4 | 7.8 | 2.5 | 23.0 |
| SPN+MCG [34] | 26.4 | 12.7 | 4.4 | 27.1 |
| PRM [33] | 44.3 | 26.8 | 9.0 | 37.6 |
| Ours | **48.5** | **30.2** | **14.4** | **44.3** |

Table 5. Image-level supervised instance segmentation results on the PASCAL VOC 2012 val. set in terms of mean average precision (mAP%) and Average Best Overlap(ABO). Our approach ourperforms the state-of-the-art PRM [33] with a relative gain of 17.8% in terms of ABO.

(PL) supervised method [14]. Our approach performs accurate counting on various categories (fruits to animals) under heavy occlusions. Fig. 6 shows counting performance comparison in terms of RMSE, across all categories, on COCO count-test. The x-axis shows different ground-truth count values. We compare with the different IC, BB and PL supervised methods [4, 14]. Our approach achieves superior results on all count values compared to glancing method [4] despite not using the beyond subitizing range annotations during training. Furthermore, we perform favourably compared to other methods using higher supervision.

**Evaluation of density map:** We employ a standard grid average mean absolute error (GAME) evaluation metric [10] used in crowd counting to evaluate spatial distribution consistency in the density map. In GAME(n), an image is divided into $4^n$ non-overlapping grid cells. Mean absolute error (MAE) between the predicted and the ground-truth local counts are reported for $n = 0, 1, 2$ and 3, as in [10]. We compare our approach with the state-of-the-art PL supervised counting approach (LCFCN) [14] on the 20 categories of the PASCAL VOC 2007 count-test set. Furthermore, we also compare with recent crowd counting approach (CSRnet) [16] on the person category of the PASCAL VOC 2007 by retraining it on the dataset. For the person category, the PL supervised LCFCN and CSRnet approaches achieve scores of 2.80 and 2.44 in GAME(3).The proposed method outperforms LCFCN and CSRnet in GAME (3) with score of 1.83, demonstrating the capabilities of our approach in the precise spatial distribution of object counts. Moreover, our method outperforms LCFCN for all 20 categories.

## 4.2. Image-level supervised Instance segmentation

Finally, we evaluate the effectiveness of our density map to improve the state-of-the-art image-level supervised instance segmentation approach (PRM) [33] on the PASCAL VOC 2012 dataset (see Sec. 3.3). In addition to PRM, the image-level supervised object detection methods MELM [27], CAM [32] and SPN [34] used with MCG mask and reported by [33] are also included in Tab. 5.

The proposed method largely outperforms all the baseline approaches and [33], in all four evaluation metrics. Even though our approach marginally increases the level of supervision (lower-count information), it improves the state-of-the-art PRM with a relative gain of 17.8% in terms of average best overlap (ABO). Compared to PRM, the gain obtained at lower IoU threshold (0.25) highlights the improved location prediction capabilities of the proposed method. Furthermore, the gain obtained at higher IoU threshold (0.75), indicates the effectiveness of the proposed scoring function in assigning higher scores to the object proposal that has highest overlap with the ground-truth object, as indicated by the improved ABO performance. Fig. 7 shows qualitative instance segmentation comparison between our approach and PRM.

## 5. Conclusion

We proposed an ILC supervised density map estimation approach for common object counting in natural scenes.

Different to existing methods, our approach provides both the global object count and the spatial distribution of object instances with the help of a novel loss function. We further demonstrated the applicability of the proposed density map in instance segmentation. Our approach outperforms existing methods for both common object counting and image-level supervised instance segmentation.

# References

[1] A. Bearman, O. Russakovsky, V. Ferrari, and L. Fei-Fei. Whats the point: Semantic segmentation with point supervision. In *ECCV*, 2016. 3

[2] S. T. Boysen and E. J. Capaldi. *The development of numerical competence: Animal and human models*. Psychology Press, 2014. 2

[3] X. Cao, Z. Wang, Y. Zhao, and F. Su. Scale aggregation network for accurate and efficient crowd counting. In *ECCV*, 2018. 1

[4] P. Chattopadhyay, R. Vedantam, R. R. Selvaraju, D. Batra, and D. Parikh. Counting everyday objects in everyday scenes. In *CVPR*, 2017. 1, 2, 3, 6, 7, 8

[5] Chhavi. k largest(or smallest) elements in an array — added min heap method, 2018. 4

[6] D. H. Clements. Subitizing: What is it? why teach it? *Teaching children mathematics*, 5, 1999. 2

[7] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *IJCV*, 111(1), 2015. 6

[8] M. Gao, A. Li, R. Yu, V. I. Morariu, and L. S. Davis. C-wsl: Count-guided weakly supervised localization. In *ECCV*, 2018. 3

[9] R. Girshick. Fast r-cnn. In *ICCV*, 2015. 3

[10] R. Guerrero, B. Torre, R. Lopez, S. Maldonado, and D. Onoro. Extremely overlapping vehicle counting. In *IbPRIA*, 2015. 8

[11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3

[12] B. R. Jansen, A. D. Hofman, M. Straatemeier, B. M. van Bers, M. E. Raijmakers, and H. L. van der Maas. The role of pattern recognition in children's exact enumeration of small numbers. *British Journal of Developmental Psychology*, 32(2), 2014. 2

[13] M. Lapin, M. Hein, and B. Schiele. Analysis and optimization of loss functions for multiclass, top-k, and multilabel classification. *TPAMI*, 40(7), 2018. 4

[14] 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. 3, 6, 7, 8

[15] V. Lempitsky and A. Zisserman. Learning to count objects in images. In *NIPS*. 2010. 3, 4

[16] Y. Li, X. Zhang, and D. Chen. Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes. In *CVPR*, 2018. 1, 8

[17] 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. 6

[18] X. Liu, J. van de Weijer, and A. D. Bagdanov. Leveraging unlabeled data for crowd counting by learning to rank. In *CVPR*, 2018. 1

[19] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 3

[20] G. Mandler and B. J. Shebo. Subitizing: an analysis of its component processes. *Journal of Experimental Psychology: General*, 111(1), 1982. 2

[21] K.-K. Maninis, J. Pont-Tuset, P. Arbeláez, and L. Van Gool. Convolutional oriented boundaries. In *ECCV*. Springer, 2016. 2, 6

[22] M. Mark, M. Kevin, L. Suzanne, and O. NoelE. Fully convolutional crowd counting on highly congested scenes. In *ICCV*, 2017. 1

[23] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Is object localization for free?-weakly-supervised learning with convolutional neural networks. In *CVPR*, 2015. 2, 4

[24] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. 5

[25] J. Pont-Tuset, P. Arbelaez, J. T. Barron, F. Marques, and J. Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *TPAMI*, 39(1), 2017. 2, 6

[26] J. Pont-Tuset and L. Van Gool. Boosting object proposals: From pascal to coco. In *ICCV*, 2015. 6

[27] F. Wan, P. Wei, J. Jiao, Z. Han, and Q. Ye. Min-entropy latent model for weakly supervised object detection. In *CVPR*, 2018. 8

[28] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu. Learning fine-grained image similarity with deep ranking. In *CVPR*, 2014. 5

[29] K. Wu, E. Otoo, and A. Shoshani. Optimizing connected component labeling algorithms. In *Medical Imaging 2005: Image Processing*, volume 5747. International Society for Optics and Photonics, 2005. 3

[30] J. Zhang, S. Ma, M. Sameki, S. Sclaroff, M. Betke, Z. Lin, X. Shen, B. Price, and R. Mech. Salient object subitizing. In *CVPR*, 2015. 3

[31] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma. Single-image crowd counting via multi-column convolutional neural network. In *CVPR*, 2016. 1

[32] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016. 4, 7, 8

[33] Y. Zhou, Y. Zhu, Q. Ye, Q. Qiu, and J. Jiao. Weakly supervised instance segmentation using class peak response. In *CVPR*, 2018. 1, 2, 4, 5, 6, 7, 8

[34] Y. Zhu, Y. Zhou, Q. Ye, Q. Qiu, and J. Jiao. Soft proposal networks for weakly supervised object localization. In *ICCV*, 2017. 8