

---

# Distribution Matching for Crowd Counting

---

Boyuan Wang\* Huidong Liu\* Dimitris Samaras Minh Hoai

Department of Computer Science, Stony Brook University, Stony Brook, NY 11790

{boywang, huidliu, samaras, minhhoai}@cs.stonybrook.edu

\*indicates equal contribution

## Abstract

In crowd counting, each training image contains multiple people, where each person is annotated by a dot. Existing crowd counting methods need to use a Gaussian to smooth each annotated dot or to estimate the likelihood of every pixel given the annotated point. In this paper, we show that imposing Gaussians to annotations hurts generalization performance. Instead, we propose to use Distribution Matching for crowd COUNTing (DM-Count). In DM-Count, we use Optimal Transport (OT) to measure the similarity between the normalized predicted density map and the normalized ground truth density map. To stabilize OT computation, we include a Total Variation loss in our model. We show that the generalization error bound of DM-Count is tighter than that of the Gaussian smoothed methods. In terms of Mean Absolute Error, DM-Count outperforms the previous state-of-the-art methods by a large margin on two large-scale counting datasets, UCF-QNRF and NWPU, and achieves the state-of-the-art results on the ShanghaiTech and UCF-CC50 datasets. DM-Count reduced the error of the state-of-the-art published result by approximately 16%.

Code is available at <https://github.com/cvlab-stonybrook/DM-Count>.

## 1 Introduction

Image-based crowd counting is an important research problem with various applications in many domains including journalism and surveillance. Current state-of-the-art methods [54, 8, 25, 55, 61, 59, 17, 48, 21, 23, 36] treat crowd counting as a density map estimation problem, where a deep neural network first produces a 2D crowd density map for a given input image and subsequently estimates the total size of the crowd by summing the density values across all spatial locations of the density map. For images of large crowds, this density map estimation approach has been shown to be more robust than the detection-then-counting approach [22, 19, 62, 12] because the former is less sensitive to occlusion and it does not need to commit to binarized decisions at an early stage.

A crucial step in the development of a density map estimation method is the training of a deep neural network that maps from an input image to the corresponding annotated density map. In all existing crowd counting datasets [15, 60, 14, 51], the annotated density map for each training image is a sparse binary mask, where each individual person is marked with a single dot on their head or forehead. The spatial extent of each person is not provided, due to the laborious effort needed for delineating the spatial extent, especially when there is too much occlusion ambiguity. Given training images with dot annotation, training the density map estimation network is equivalent to optimizing the parameters of the network to minimize a differentiable loss function that measures the discrepancy between the predicted density map and the dot-annotation map. Notably, the former is a dense real-value matrix, while the later is a sparse binary matrix. Given the sparsity of the dots, a function that is defined based on the pixel-wise difference between the annotated and predicted density maps is hard to train because the reconstruction loss is heavily unbalanced between the 0s and 1s in the sparse binary matrix. One approach to alleviate this problem is to turn each annotated dot into a Gaussian blob such that the ground truth is more balanced and thus the network is easier to train. Almost all prior crowd density map estimation methods [56, 57, 60, 38, 20, 33, 35, 4, 28, 50, 27, 40, 29, 26] have followed

this convention. Unfortunately, the performance of the resulting network is highly dependent on the quality of this “pseudo ground truth”, but it is not trivial to set the right widths for the Gaussian blobs given huge variation in the sizes and shapes of people in a perspective image of a crowded scene.

Recently, Ma et al. [31] proposed a Bayesian loss to measure the discrepancy between the predicted and the annotated density maps. This method transforms a binary ground truth annotation map into  $N$  “smoothed ground truth” density maps, where  $N$  is the count number. Each pixel value of a smoothed ground truth density map is the posterior probability of the corresponding annotation dot given the location of that pixel. Empirically, this method has been shown to outperform other aforementioned approaches [60, 38, 20, 33, 35, 4]. However, there are two major problems with this loss function. First, it also requires a Gaussian kernel to construct the likelihood function for each annotated dot, which involves setting the kernel width. Second, this loss corresponds to an underdetermined system of equations with infinitely many solutions. The loss can be 0 for many density maps that are not similar to the ground truth density map. As a consequence, using this loss for training can lead to a predicted density map that is very different from the ground truth density map.

In this paper, we address the shortcomings in existing approaches with the following contributions.

- We theoretically and empirically show that imposing Gaussians to annotations will hurt the generalization performance of a crowd counting network.
- We propose DM-Count, a method that performs Distribution Matching for crowd COUNTing. Unlike previous works, DM-Count does not need any Gaussian smoothing ground truth annotations. Instead, we use Optimal Transport (OT) to measure the similarity between the normalized predicted density map and the normalized ground truth density map. To stabilize the OT computation, we further add a Total Variation (TV) loss.
- We present the generalization error bounds for the counting loss, OT loss, TV loss and the overall loss in our method. All the bounds are tighter than those of the Gaussian smoothed methods.
- Empirically, our method improved the state-of-the-art by a large margin on four challenging crowd counting datasets: UCF-QNRF, NWPU, ShanghaiTech, and UCF-CC50. Notably, our method reduced the published state-of-the-art MAE on the NWPU dataset by approximately 16%.

## 2 Previous Work

### 2.1 Crowd Counting Methods

Crowd counting methods can be divided into three categories: detection-then-count, direct count regression, and density map estimation. Early methods [22, 19, 62, 12] detect people, heads, or upper bodies in the image. However, accurate detection is difficult for dense crowds. Besides, it also requires bounding box annotation, which is a laborious and ambiguous process due to heavy occlusion. Later methods [5, 6, 49, 7] avoid the detection problem and directly learn to regress the count from a feature vector. But their results are less interpretable and the dot annotation maps are underutilized. Most recent works [20, 35, 15, 4, 31, 50, 27, 40, 29, 26, 54, 8, 25, 55, 61, 59, 17, 48, 21, 23, 30, 47, 53, 43, 37, 56, 18, 39, 24, 44, 42] are based on density map estimation, which has been shown to be more robust than detection-then-count and count regression approaches.

Density map estimation methods usually define the training loss based on the pixel-wise difference between the Gaussian smoothed density map and the predicted density map. Instead of using a single kernel width to smooth the dot annotation, [60, 14, 47] used adaptive kernel width. The kernel width is selected based on the distance to an annotated dot’s nearest neighbors. Specifically, [15] generated multiple smoothed ground truth density maps on different density levels. The final loss combines the reconstruction errors from multiple density levels. However, these methods assume the crowd is evenly distributed; in reality crowd distribution is quite irregular. The Bayesian loss method [31] uses a Gaussian to construct a likelihood function for each annotated dot. However, it may not predict a correct density because the loss is underdetermined. Detailed analysis can be found in Sec 4.2.

### 2.2 Optimal Transport

We propose a novel loss function based on Optimal Transport (OT) [46]. For a better understanding of the proposed method, we briefly review the Monge-Kantorovich OT formulation in this section.

Optimal Transport refers to the optimal cost to transform one probability distribution to another. Let  $\mathcal{X} = \{\mathbf{x}_i | \mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$  and  $\mathcal{Y} = \{\mathbf{y}_j | \mathbf{y}_j \in \mathbb{R}^d\}_{j=1}^n$  be two sets of points on  $d$ -dimensional vector space. Let  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  be two probability measures defined on  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively;  $\boldsymbol{\mu}, \boldsymbol{\nu} \in \mathbb{R}_+^n$  and  $\mathbf{1}_n^T \boldsymbol{\mu} = \mathbf{1}_n^T \boldsymbol{\nu} = 1$  ( $\mathbf{1}_n$  is a  $n$ -dimensional vector of all ones). Let  $c : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}_+$  be the cost function for moving from a point in  $\mathcal{X}$  to a point in  $\mathcal{Y}$ , and  $\mathbf{C}$  be the corresponding  $n \times n$  cost matrix for the two sets of points:  $\mathbf{C}_{ij} = c(\mathbf{x}_i, \mathbf{y}_j)$ . Let  $\Gamma$  be the set of all possible ways to transport probability mass from  $\mathcal{X}$  to  $\mathcal{Y}$ :  $\Gamma = \{\gamma \in \mathbb{R}_+^{n \times n} : \gamma \mathbf{1} = \boldsymbol{\mu}, \gamma^T \mathbf{1} = \boldsymbol{\nu}\}$ . The Monge-Kantorovich’s Optimal Transport (OT) cost between  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  is defined as:

$$\mathcal{W}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\gamma \in \Gamma} \langle \mathbf{C}, \gamma \rangle. \quad (1)$$

Intuitively, if the probability distribution  $\boldsymbol{\mu}$  is viewed as a unit amount of “dirt” piled on  $\mathcal{X}$  and  $\boldsymbol{\nu}$  a unit amount of dirt piled on  $\mathcal{Y}$ , the OT cost is the minimum “cost” of turning one pile into the other. The OT cost is a principal measurement to quantify the dissimilarity between two probability distributions, also taking into account the distance between “dirt” locations.

The OT cost can also be computed via the dual formulation:

$$\mathcal{W}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{R}^n} \langle \boldsymbol{\alpha}, \boldsymbol{\mu} \rangle + \langle \boldsymbol{\beta}, \boldsymbol{\nu} \rangle, \quad \text{s.t. } \alpha_i + \beta_j \leq c(\mathbf{x}_i, \mathbf{y}_j), \quad \forall i, j. \quad (2)$$

### 3 DM-Count: Distribution Matching for Crowd Counting

We consider crowd counting as a distribution matching problem. In this section, we propose DM-Count: Distribution matching for crowd counting. A network for crowd counting inputs an image and outputs a map of density values. The final count estimate can be obtained by summing over the predicted density map. DM-Count is agnostic to different network architectures. In our experiments, we use the same network as in the Bayesian loss paper [31]. Unlike all previous density map estimation methods which need to use Gaussians to smooth ground truth annotations, DM-Count does not need any Gaussian to preprocess ground truth annotations.

Let  $\mathbf{z} \in \mathbb{R}_+^n$  denote the vectorized binary map for dot-annotation and  $\hat{\mathbf{z}} \in \mathbb{R}_+^n$  the vectorized predicted density map returned by a neural network. By viewing  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  as unnormalized density functions, we formulate the loss function in DM-Count using three terms: the counting loss, the OT loss, and the Total Variation (TV) loss. The first term measures the difference between the total masses, while the last two measures the difference between the distributions of the normalized density functions.

**The Counting Loss.** Let  $\|\cdot\|_1$  denote the  $L_1$  norm of a vector, and so  $\|\mathbf{z}\|_1, \|\hat{\mathbf{z}}\|_1$  are the ground truth and predicted counts respectively. The goal of crowd counting is to make  $\|\hat{\mathbf{z}}\|_1$  as close as possible to  $\|\mathbf{z}\|_1$ , and the counting loss is defined as the absolute difference between them:

$$\ell_C(\mathbf{z}, \hat{\mathbf{z}}) = \left| \|\mathbf{z}\|_1 - \|\hat{\mathbf{z}}\|_1 \right|. \quad (3)$$

**The Optimal Transport Loss.** Both  $\mathbf{z}$  and  $\hat{\mathbf{z}}$  are unnormalized density functions, but we can turn them into probability density functions (pdfs) by dividing them by their respective total mass. Apart from OT, the Kullback-Leibler divergence and Jensen-Shannon divergence can also measure the similarity between two pdfs. However, these measurements do not provide valid gradients to train a network if the source distribution does not overlap with the target distribution [32]. Therefore, we propose the use of OT in this work. We define the OT loss as follows:

$$\ell_{OT}(\mathbf{z}, \hat{\mathbf{z}}) = \mathcal{W}\left(\frac{\mathbf{z}}{\|\mathbf{z}\|_1}, \frac{\hat{\mathbf{z}}}{\|\hat{\mathbf{z}}\|_1}\right) = \left\langle \boldsymbol{\alpha}^*, \frac{\mathbf{z}}{\|\mathbf{z}\|_1} \right\rangle + \left\langle \boldsymbol{\beta}^*, \frac{\hat{\mathbf{z}}}{\|\hat{\mathbf{z}}\|_1} \right\rangle, \quad (4)$$

where  $\boldsymbol{\alpha}^*$  and  $\boldsymbol{\beta}^*$  are the solutions of Problem (2). We use the quadratic transport cost, i.e.,  $c(\mathbf{z}(i), \hat{\mathbf{z}}(j)) = \|\mathbf{z}(i) - \hat{\mathbf{z}}(j)\|_2^2$ , where  $\mathbf{z}(i)$  and  $\hat{\mathbf{z}}(j)$  are 2D coordinates of locations  $i$  and  $j$ , respectively. To avoid the division-by-zero error, we add a machine precision to the denominator.

Since the entries in  $\hat{\mathbf{z}}$  are non-negative, the gradient of Eq. (4) with respect to  $\hat{\mathbf{z}}$  is:

$$\frac{\partial \ell_{OT}(\mathbf{z}, \hat{\mathbf{z}})}{\partial \hat{\mathbf{z}}} = \frac{\boldsymbol{\beta}^*}{\|\hat{\mathbf{z}}\|_1} - \frac{\langle \boldsymbol{\beta}^*, \hat{\mathbf{z}} \rangle}{\|\hat{\mathbf{z}}\|_1^2}. \quad (5)$$

This gradient can be back-propagated to learn the parameters of the density estimation network.

**Total Variation Loss.** In each training iteration, we use the Sinkhorn algorithm [34] to approximate  $\alpha^*$  and  $\beta^*$ . The time complexity is  $O(n^2 \log n / \epsilon^2)$  [9], where  $\epsilon$  is the desired optimality gap, i.e., the upper bound for the difference between the returned objective and the optimal objective. When optimizing with the Sinkhorn algorithm, the objective decreases dramatically at the beginning but only converges slowly to the optimal objective in later iterations. In practice, we set the maximum number of iterations, and the Sinkhorn algorithm only returns an approximate solution. As a result, when we optimize the OT loss with the Sinkhorn algorithm, the predicted density map ends up close to the ground truth density map, but not exactly the same. The OT loss will approximate well the dense areas of the crowd, but the approximation might be poorer for the low density areas of the crowd. To address this issue, we additionally use the Total Variation (TV) loss, defined as<sup>1</sup>:

$$\ell_{TV}(\mathbf{z}, \hat{\mathbf{z}}) = \left\| \frac{\mathbf{z}}{\|\mathbf{z}\|_1} - \frac{\hat{\mathbf{z}}}{\|\hat{\mathbf{z}}\|_1} \right\|_{TV} = \frac{1}{2} \left\| \frac{\mathbf{z}}{\|\mathbf{z}\|_1} - \frac{\hat{\mathbf{z}}}{\|\hat{\mathbf{z}}\|_1} \right\|_1. \quad (6)$$

The TV loss will also increase the stability of the training procedure. Optimizing the OT loss with the Sinkhorn algorithm is a min-max saddle point optimization procedure, which is similar to GAN optimization [13]. The stability of GAN training can be increased by adding a reconstruction loss, as shown in the Pix2Pix GAN [16]. To this end, the TV loss is similar to the reconstruction loss, and also increases the stability of the training procedure.

The gradient of the TV loss with respect to the predicted density map  $\hat{\mathbf{z}}$  is:

$$\frac{\partial \ell_{TV}(\mathbf{z}, \hat{\mathbf{z}})}{\partial \hat{\mathbf{z}}} = -\frac{1}{2} \left( \frac{\text{sign}(\mathbf{v})}{\|\hat{\mathbf{z}}\|_1} - \frac{\langle \text{sign}(\mathbf{v}), \hat{\mathbf{z}} \rangle}{\|\hat{\mathbf{z}}\|_1^2} \right), \quad (7)$$

where  $\mathbf{v} = \mathbf{z} / \|\mathbf{z}\|_1 - \hat{\mathbf{z}} / \|\hat{\mathbf{z}}\|_1$ , and  $\text{sign}(\cdot)$  is the Sign function on each element of a vector.

**The Overall Objective.** The overall loss function is the combination of the counting loss, the OT loss, and the TV loss:

$$\ell(\mathbf{z}, \hat{\mathbf{z}}) = \ell_C(\mathbf{z}, \hat{\mathbf{z}}) + \lambda_1 \ell_{OT}(\mathbf{z}, \hat{\mathbf{z}}) + \lambda_2 \|\mathbf{z}\|_1 \ell_{TV}(\mathbf{z}, \hat{\mathbf{z}}), \quad (8)$$

where  $\lambda_1$  and  $\lambda_2$  are tunable hyper-parameters for the OT and TV losses. To ensure that the TV loss has the same scale as the counting loss, we multiply this loss term with the total count.

Given  $K$  training images  $\{I_k\}_{k=1}^K$  with corresponding dot annotation maps  $\{\mathbf{z}_k\}_{k=1}^K$ , we will learn a deep neural network  $f$  for density map estimation by minimizing:  $L(f) = \frac{1}{K} \sum_{k=1}^K \ell(\mathbf{z}_k, f(I_k))$ .

## 4 Generalization Bounds and Theoretical Analysis

In this section, we analyze the theoretical properties of the Gaussian smoothed methods, the Bayesian loss, and the proposed DM-Count. The proofs of the theorems in this section can be found in the supplementary material. First, we introduce some notations below.

Let  $\mathcal{I}$  denote the set of images and  $\mathcal{Z}$  the set of dot annotation maps. Let  $\mathcal{D} = \{(I, \mathbf{z})\}$  be the joint distribution of crowd images and corresponding dot annotation maps. Let  $\mathcal{H}$  be a hypothesis space. Each  $h \in \mathcal{H}$  maps from  $I \in \mathcal{I}$  to each dimension of  $\mathbf{z} \in \mathcal{Z}$ . Let  $\mathcal{F} = \mathcal{H} \times \dots \times \mathcal{H}$  ( $n$  times) be the mapping space. Each  $f \in \mathcal{F}$  maps  $I \in \mathcal{I}$  to  $\mathbf{z} \in \mathcal{Z}$ . Let  $\mathbf{t}$  be the Gaussian smoothed density map of each  $\mathbf{z} \in \mathcal{D}$ , and let  $\tilde{\mathcal{D}} = \{(I, \mathbf{t})\}$  be the joint distribution of  $(I, \mathbf{t})$ . Let  $S = \{(I_k, \mathbf{z}_k)\}_{k=1}^K$ , and  $\tilde{S} = \{(I_k, \mathbf{t}_k)\}_{k=1}^K$  be the finite sets of  $K$  samples i.i.d. sampled from  $\mathcal{D}$  and  $\tilde{\mathcal{D}}$ , respectively. Let  $R_S(\mathcal{H})$  denote the empirical Rademacher complexity [3] for  $\mathcal{H}$  w.r.t  $S$ . Given a data set  $D \in \{\mathcal{D}, S, \tilde{\mathcal{D}}, \tilde{S}\}$ , a mapping  $f \in \mathcal{F}$  and a loss function  $\ell$ , let  $\mathcal{R}(D, f, \ell) = \mathbb{E}_{(I, \mathbf{s}) \sim D}[\ell(\mathbf{s}, f(I))]$  denote the expected risk. Let  $\ell_1(\mathbf{z}, \hat{\mathbf{z}}) = \|\mathbf{z} - \hat{\mathbf{z}}\|_1$ . Let  $f_{\Delta}^D = \text{argmin}_{f \in \mathcal{F}} \mathcal{R}(D, f, \ell_{\Delta})$  be the minimizer of  $\mathcal{R}(D, f, \ell_{\Delta})$  over a data set  $D$  using the loss  $\ell_{\Delta}$ , where  $D \in \{\mathcal{D}, S, \tilde{\mathcal{D}}, \tilde{S}\}$ , and  $\Delta \in \{1, C, OT, TV, \emptyset\}$ .

### 4.1 Generalization Error Bounds of Gaussian Smoothed Methods

Many existing methods (e.g., [60, 20, 35]) use Gaussian-smoothed annotation maps for training. Below we give generalization error bounds when using the  $\ell_1$  loss on the density maps.

<sup>1</sup>In the training loss context, Total Variation refers to the total variation distance of two probability measures. A formal definition can be found in [45, Definition 2.4, page 83]. Eq. (6) is [45, Lemma 2.1, page 84].

**Theorem 1** Assume that  $\forall f \in \mathcal{F}$  and  $(I, \mathbf{t}) \sim \tilde{\mathcal{D}}$ , we have  $\ell(\mathbf{t}, f(I)) \leq B$ . Then, for any  $0 < \delta < 1$ , with probability of at least  $1 - \delta$ ,

a) the upper bound of the generalization error is

$$\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1) \leq \mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{D}}, \ell_1) + 2nR_{\tilde{S}}(\mathcal{H}) + 5B\sqrt{2\log(8/\delta)/K} + \mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1,$$

b) the lower bound of the generalization error is

$$\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1) \geq \left| \mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1 - \mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1) \right|.$$

In this theorem, as the number of samples  $K$  grows to infinity,  $2nR_{\tilde{S}}(\mathcal{H})$  and  $5B\sqrt{2\log(8/\delta)/K}$  decrease to 0. Theorem 1.a) shows that the upper bound (worst case) of the expected risk  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1)$ , which is evaluated on real ground truth data using an empirical minimizer trained on the Gaussian smoothed ground truth, does not exceed  $\mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{D}}, \ell_1) + \mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1$  given sufficient training data. Theorem 1.b) shows that the lower bound (best case) of  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1)$  is not smaller than  $|\mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1 - \mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1)|$ . This means that if  $\mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1) \leq \mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1$ , then the smaller  $\mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1)$  is, the larger the expected risk  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1)$  will be. In other words, the better a good model  $f_1^{\tilde{S}}$  performs on the Gaussian smoothed ground truth  $\tilde{\mathcal{D}}$ , the poorer it generalizes on the real ground truth  $\mathcal{D}$ . Furthermore, as long as  $\mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1) \neq \mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1$ , we have  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1) > 0$ .  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1)$  can be as large as  $\mathbb{E}_{(I, \mathbf{z}) \sim \mathcal{D}} \|\mathbf{z} - \mathbf{t}\|_1$  when  $\mathcal{R}(\tilde{\mathcal{D}}, f_1^{\tilde{S}}, \ell_1) = 0$ . This is undesirable because we want the risk  $\mathcal{R}(\mathcal{D}, f_1^{\tilde{S}}, \ell_1)$  evaluated on the real ground truth to be 0 as well.

## 4.2 The Underdetermined Bayesian Loss

The Bayesian Loss [31] is:

$$\ell_{\text{Bayesian}}(\mathbf{z}, \hat{\mathbf{z}}) = \sum_{i=1}^N |1 - \langle \mathbf{p}_i, \hat{\mathbf{z}} \rangle|, \text{ where } \mathbf{p}_i = \frac{\mathcal{N}(\mathbf{q}_i, \sigma^2 \mathbf{1}_{2 \times 2})}{\sum_{i=1}^N \mathcal{N}(\mathbf{q}_i, \sigma^2 \mathbf{1}_{2 \times 2})}, \quad (9)$$

and  $N$  is number of people of  $\mathbf{z}$ , and  $\mathcal{N}(\mathbf{q}_i, \sigma^2 \mathbf{1}_{2 \times 2})$  is a Gaussian distribution centered at  $\mathbf{q}_i$  with variance  $\sigma^2 \mathbf{1}_{2 \times 2}$ .  $\mathbf{q}_i$  is the  $i^{\text{th}}$  annotated dot in  $\mathbf{z}$ . The dimension of  $\mathbf{p}_i$  and  $\mathbf{z}$  is  $n$ , the number of pixels of the density map. However, since the number of annotated dots  $N$  is less than  $n$ , the Bayesian loss is underdetermined. For a ground truth annotation  $\mathbf{z}$ , there are infinitely many  $\hat{\mathbf{z}}$  with  $\ell_{\text{Bayesian}}(\mathbf{z}, \hat{\mathbf{z}}) = 0$  and  $\hat{\mathbf{z}} \neq \mathbf{z}$ . Therefore, the predicted density map could be very different from the ground truth density map.

## 4.3 The Generalization Error Bounds of the Losses in DM-Count

We give the generalization error bounds of the losses in the proposed method in the following theorem.

**Theorem 2** Assume that  $\forall f \in \mathcal{F}$  and  $(I, \mathbf{z}) \sim \mathcal{D}$ , we have  $\|\mathbf{z}\|_1 \geq 1$ ,  $\|f(I)\|_1 \geq 1$  (can be satisfied by adding a dummy dimension with value of 1 to both  $\mathbf{z}$  and  $f(I)$ ) and  $\ell_C(\mathbf{z}, f(I)) \leq B$ . Then, for any  $0 < \delta < 1$ , with probability of at least  $1 - \delta$

a) the generalization error bound of the counting loss is

$$\mathcal{R}(\mathcal{D}, f_C^S, \ell_C) \leq \mathcal{R}(\mathcal{D}, f_C^D, \ell_C) + 2nR_S(\mathcal{H}) + 5B\sqrt{2\log(8/\delta)/K},$$

b) the generalization error bound of the OT loss is

$$\mathcal{R}(\mathcal{D}, f_{OT}^S, \ell_{OT}) \leq \mathcal{R}(\mathcal{D}, f_{OT}^D, \ell_{OT}) + 4\mathbf{C}_{\infty} n^2 R_S(\mathcal{H}) + 5\mathbf{C}_{\infty} \sqrt{2\log(8/\delta)/K},$$

c) the generalization error bound of the TV loss is

$$\mathcal{R}(\mathcal{D}, f_{TV}^S, \ell_{TV}) \leq \mathcal{R}(\mathcal{D}, f_{TV}^D, \ell_{TV}) + n^2 R_S(\mathcal{H}) + 5\sqrt{2\log(8/\delta)K},$$

d) the generalization error bound of the overall loss is

$$\begin{aligned} \mathcal{R}(\mathcal{D}, f^S, \ell) &\leq \mathcal{R}(\mathcal{D}, f^D, \ell) + (2n + 4\lambda_1 \mathbf{C}_{\infty} n^2 + \lambda_2 N n^2) R_S(\mathcal{H}) \\ &\quad + 5(B + \lambda_1 \mathbf{C}_{\infty} + \lambda_2 N) \sqrt{2\log(8/\delta)K}, \end{aligned}$$

where  $\mathbf{C}_{\infty}$  is the maximum cost in the cost matrix in OT, and  $N = \sup\{\|\mathbf{z}\|_1 \mid \forall (I, \mathbf{z}) \sim \mathcal{D}\}$  is the maximum count number over a dataset.

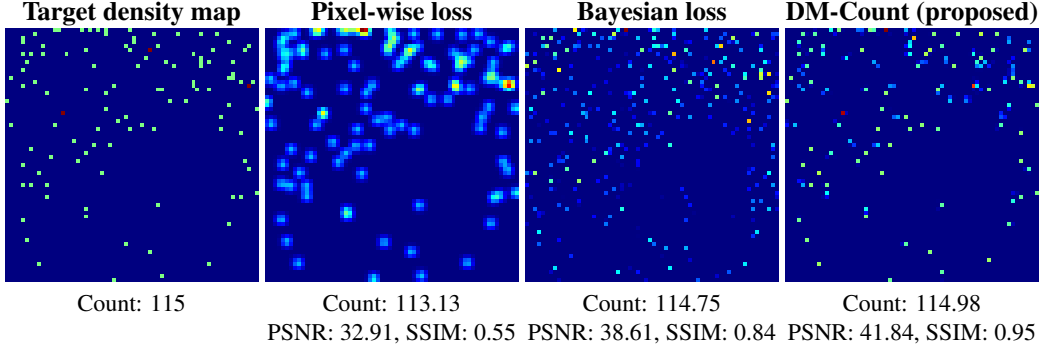


Figure 1: **Comparison of different methods on toy data.** The pixel-wise loss generates a blurry density map with a higher counting error. The Bayesian loss produces dissimilar density maps from the ground truth, with high values in many locations with no annotations. DM-Count is able to produce more accurate crowd count and localization than the other two methods.

In the above theorem, as  $K$  grows,  $R_S(\mathcal{H})$  and  $\sqrt{2\log(1/\delta)K}$  decrease. All the expected risks  $\mathcal{R}(\mathcal{D}, f_{\Delta}^S, \ell_{\Delta})$  using the empirical minimizers  $f_{\Delta}^S$  converge to the expected risks  $\mathcal{R}(\mathcal{D}, f_{\Delta}^D, \ell_{\Delta})$ ,  $\Delta \in \{C, OT, TV, \emptyset\}$  using optimal minimizers  $f_{\Delta}^D$ . This means that all the upper bounds are tight. In addition, all upper bounds are tighter than the upper bound of the Gaussian smoothed methods shown in Theorem 1.a). The bound of the OT loss in Theorem 2.b) is related to the maximum transport cost  $C_{\infty}$ . Therefore, we need to use a smaller transport cost in OT for better generalization performance. The coefficient of  $R_S(\mathcal{H})$  for the counting loss is  $O(n)$ , and for the OT loss and the TV loss is  $O(n^2)$ . This means that for larger image size, we need more images to train. The number is linear to the size of  $\mathbf{z}$  using solely the counting loss, and quadratic using solely the OT loss or the TV loss. When using all three losses, we need to set  $\lambda_1$  and  $\lambda_2$  to be small in order to balance the three losses.

## 5 Experiments

In this section, we describe experiments on toy data and on benchmark crowd counting datasets. More detailed dataset descriptions, implementation details and experimental settings can be found in the supplementary material.

### 5.1 Results on Toy Data

To understand the empirical behavior of different methods, we consider a toy problem where the task is to move a source density map  $\hat{\mathbf{z}}$  to a target density map  $\mathbf{z}$  using the Pixel-wise loss, the Bayesian loss and DM-Count. The source density map  $\hat{\mathbf{z}}$  is initialized from a uniform distribution between 0 and 0.01, and the target density map is shown in the leftmost figure in Fig. 1. All three methods start from the same source density map. Fig. 1 visualizes the final  $\hat{\mathbf{z}}$  at convergence. The Pixel-wise loss yields a blurry density map with a higher count. The Bayesian loss performs better than the Pixel-wise loss in terms of counting error, Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity in Image (SSIM) [52], but the resulting density map is quite different from the target, with high values at many locations where no dots are annotated. This confirms our analysis that the Bayesian loss corresponds to an underdetermined system such that the output density map could be very different from the target density map. In contrast, DM-Count is able to produce a more accurate count and density map. DM-Count outperforms the Bayesian loss by a large margin in both PSNR and SSIM.

### 5.2 Results on Benchmark Datasets

We perform experiments on four challenging crowd counting datasets: UCF-QNRF [15], NWPU [51], ShanghaiTech [60], and UCF-CC-50 [14]. It is worth noting that the NWPU dataset is the largest-scale and most challenging crowd counting dataset publicly available today. The ground truth counts for test images are not released, and the results on the test set must be obtained by submitting to the evaluation server at <https://www.crowdbenchmark.com/nwpucrowd.html>. Following previous work [35, 15, 4, 14, 60], we use the following metrics: Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and Mean Normalized Absolute Error (NAE) as evaluation metrics. For

	UCF-QNRF		ShanghaiTech A		ShanghaiTech B		UCF-CC-50	
	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE
Crowd CNN [58]	-	-	181.8	277.7	32.0	49.8	467.0	498.5
MCNN [60]	277	426	110.2	173.2	26.4	41.3	377.6	509.1
CMTL [41]	252	514	101.3	152.4	20.0	31.1	322.8	341.4
Switch CNN [2]	228	445	90.4	135.0	21.6	33.4	318.1	439.2
IG-CNN [1]	-	-	72.5	118.2	13.6	21.1	291.4	349.4
ic-CNN [35]	-	-	68.5	116.2	10.7	16.0	260.9	365.5
CSR Net [20]	-	-	68.2	115.0	10.6	16.0	266.1	397.5
SANet [4]	-	-	67.0	104.5	8.4	13.6	258.4	334.9
CL-CNN [15]	132	191	-	-	-	-	-	-
PACNN [40]	-	-	62.4	102.0	7.6	11.8	241.7	320.7
CAN [27]	107	183	62.3	100.0	7.8	12.2	212.2	<b>243.7</b>
SFCN [50]	102	171	64.8	107.5	7.6	13.0	214.2	318.2
ANF [57]	110	174	63.9	99.4	8.3	13.2	250.2	340.0
Wan <i>et al.</i> [47]	101	176	64.7	97.1	8.1	13.6	-	-
Pixel-wise Loss [31]	106.8	183.7	68.6	110.1	8.5	13.9	251.6	331.3
Bayesian Loss [31]	88.7	154.8	62.8	101.8	7.7	12.7	229.3	308.2
DM-Count (proposed)	<b>85.6</b>	<b>148.3</b>	<b>59.7</b>	<b>95.7</b>	<b>7.4</b>	<b>11.8</b>	<b>211.0</b>	291.5

Table 1: Results on the UCF-QNRF, Shanghai Tech, and UCF-CC-50 datasets.

	Backbone	Validation set		Test set		
		MAE	RMSE	MAE	RMSE	NAE
MCNN [60]	FS	218.5	700.6	232.5	714.6	1.063
CSR net [20]	VGG-16	104.8	433.4	121.3	387.8	0.604
PCC-Net-VGG [10]	VGG-16	100.7	573.1	112.3	457.0	0.251
CAN [27]	VGG-16	93.5	489.9	106.3	<b>386.5</b>	0.295
SCAR [11]	VGG-16	81.5	397.9	110.0	495.3	0.288
Bayesian Loss [31]	VGG-19	93.6	470.3	105.4	454.2	0.203
SFCN [50]	ResNet-101	95.4	608.3	105.7	424.1	0.254
DM-Count (proposed)	VGG-19	<b>70.5</b>	<b>357.6</b>	<b>88.4</b>	388.6	<b>0.169</b>

Table 2: Results of various methods on the NWPU validation and test sets.

all three metrics, the smaller the better. For a fair comparison, we use the same network as in the Bayesian loss paper [31]. In all experiments, we set  $\lambda_1 = 0.1$ ,  $\lambda_2 = 0.01$ , and the Sinkhorn entropic regularization parameter to 10. The number of Sinkhorn iterations is set to 100. On average, the OT computation time is  $25ms$  for each image.

**Quantitative Results.** Tables 1 and 2 compare the performance of DM-Count against various methods. In all experiments, DM-Count outperforms all other methods except CAN under MSE in NWPU (where they are comparable). Although we use the same set of hyper-parameters for DM-Count in all experiments, DM-Count still achieves the best performance, suggesting that DM-Count’s performance is stable across various datasets.

DM-Count outperforms the Pixel-wise loss and the Bayesian loss, when used in the same network architecture and training procedure as DM-Count, in all the experiments. This demonstrates the effectiveness of the proposed loss. The pixel-wise loss is much worse than DM-Count in Table 1. Additionally, even without using a multi-scale architecture as in [4, 47], or a deeper network as in [2, 50], DM-Count still achieves state-of-the-art performance on all four datasets. This indicates the importance of having a good loss function in crowd counting.

On the large-scale and challenging datasets UCF-QNRF and NWPU, DM-Count significantly outperforms the state-of-the-art methods. Specifically, on the UCF-QNRF dataset, DM-Count reduces the MAE and MSE of the Bayesian loss from 88.7 to 85.6 and from 154.8 to 148.3, respectively. Notably, on the NWPU test set (obtained by submitting to the evaluation server), DM-Count reduces the MAE and NAE by a large margin, from 105.4 to 88.4 in MAE and from 0.203 to 0.169 in NAE.

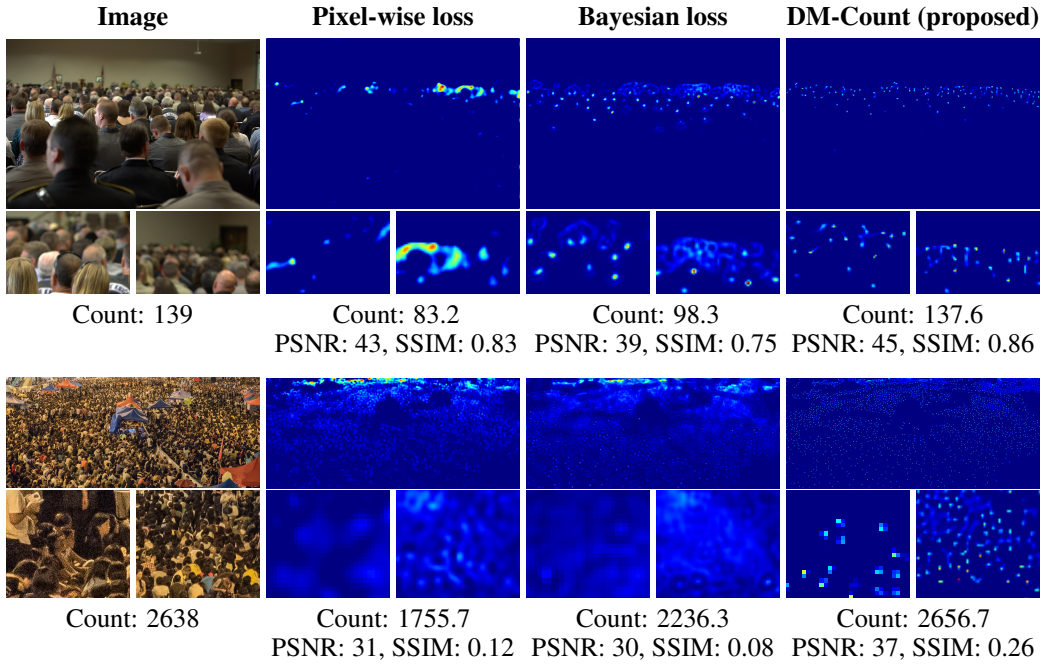


Figure 2: **Density map visualization.** Comparison between Pixel-wise loss, Bayesian loss and DM-Count. The pixel-wise and Bayesian losses fail to localize people well in dense regions. DM-Count is able to localize people both in dense and sparse regions. The Count number, PSNR and SSIM metrics suggest that DM-Count produces more accurate count numbers and better density maps.

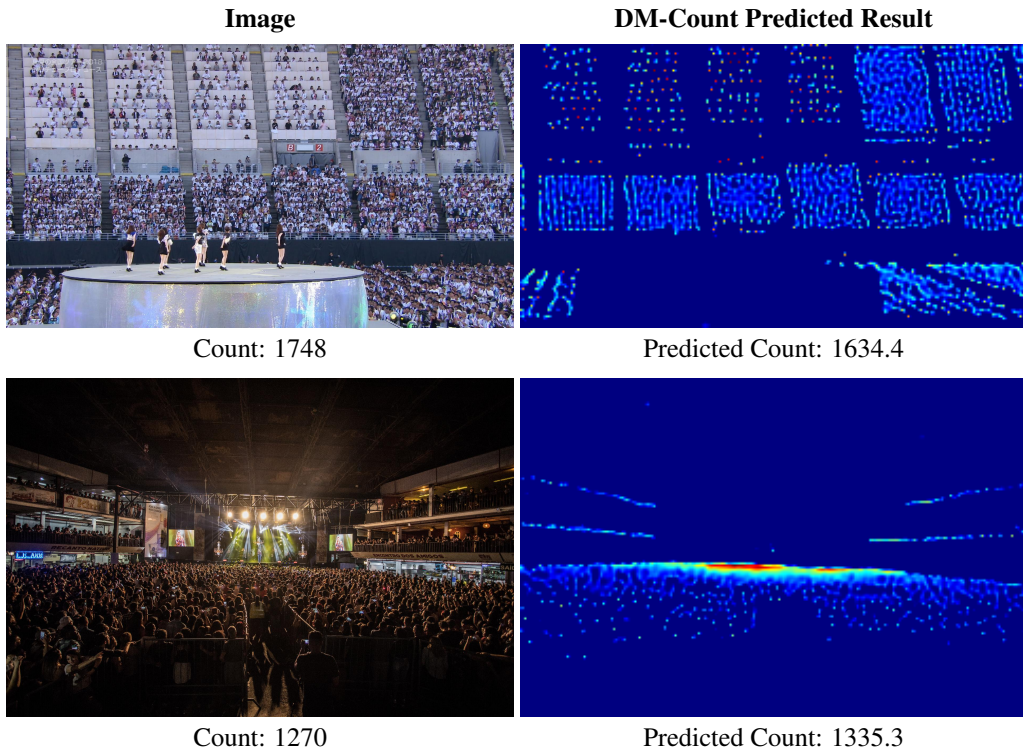


Figure 3: **Density map visualization on the NWPU validation set.**

**Qualitative Results.** Fig. 2 shows the predicted density maps of the Pixel-wise loss, the Bayesian loss and DM-Count. This figure demonstrates that: 1) DM-Count produces count numbers that are closer to the ground truth numbers, 2) DM-Count produces much sharper density maps than the



# Sinkhorn Iters	MAE	RMSE
50	90.8	162.1
100	85.6	148.3
120	85.5	151.5

Table 3: **Effect of # of Sinkhorn iterations.**

Method	MAE	RMSE
Pixel-wise loss	144.1	232.5
Bayesian loss	108.4	187.2
DM-Count	<b>105.6</b>	<b>181.6</b>

Table 4: **Robustness to noisy annotations.**

Pixel-wise and Bayesian losses. In Fig. 2, DM-Count produces much higher PSNRs and SSIMs than the Pixel-wise and Bayesian losses. The average PSNR and SSIM over the whole UCF-QNRF test set for the Pixel-wise loss are 34.79 and 0.43, for the Bayesian loss are 34.55 and 0.42, and for DM-Count are 40.65 and 0.55, respectively. Because the Pixel-wise loss uses the Gaussian smoothed ground truth, it produces blurrier density maps than the real ground truth. This empirically verifies our theoretical analysis of the generalization bound of Gaussian smoothed methods. As shown in the figure, the Pixel-wise and Bayesian losses are unable to localize people in dense regions. In contrast, DM-Count localizes people well in both dense and sparse regions. Fig. 3 shows predicted density maps by DM-Count. The predicted density maps correspond well to crowd densities in both sparse and dense areas, demonstrating the effectiveness of DM-Count in spatial density estimation.

### 5.3 Ablation Studies

**Hyper-parameter study.** We tune  $\lambda_1$  and  $\lambda_2$  in DM-Count on the UCF-QNRF dataset. First, we fix  $\lambda_1$  to 0.1 and tune  $\lambda_2$  from 0.01, 0.05 to 0.1. The MAE varies from 85.6, 87.8 to 88.5. As  $\lambda_2 = 0.01$  achieves the best result, we fix  $\lambda_2$  to 0.01 and tune  $\lambda_1$  from 0.01, 0.05 to 0.1. The MAE varies from 87.2, 86.2 to 85.6. Thus, we set  $\lambda_1 = 0.1$ ,  $\lambda_2 = 0.01$  and use them on all the datasets.

**Effect of the number of Sinkhorn iterations.** Table 3 lists the results of DM-Count on the UCF-QNRF dataset using different numbers of Sinkhorn iterations. As shown in this table, using a small number of iterations lowers the performance of DM-Count, which indicates that we obtain inaccurate OT solutions. When the number of iterations increases to 100, DM-Count outperforms the previous state-of-the-art. The performance plateaued after the number of iterations crossed 100. Therefore, in all of our experiments, we use 100 Sinkhorn iterations for DM-Count.

**Contribution of each component.** The loss in DM-Count is composed of three components, the counting loss, the OT loss and the TV loss. We study the contribution of each component on the UCF-QNRF dataset. Results are listed in Table 5. As seen in the Table, all components are essential to the final performance. However, the OT loss is the most important component.

Component	Combinations			
Counting loss	✓	✓	✓	✓
OT loss			✓	✓
TV loss		✓		✓
MAE	103.1	94.9	89.3	85.6
RMSE	175.9	167.4	161.3	148.3

Table 5: **Component analysis**

**Robustness to noisy annotations.** Crowd annotation is performed by placing a single dot on a person. Such process is ambiguous and could lead to inevitable annotation errors. We study how different loss functions perform w.r.t. annotation errors. We add uniform random noise to the original annotation and train different models with the same noisy annotation. The noise is randomly generated between 0 and 5% of the image height, and is about 80 pixels on average. As shown in Table 4, the proposed DM-Count is more robust to annotation errors compared to the pixel-wise Bayesian losses.

## 6 Conclusion

In this paper, we have shown that using the Gaussian kernel to smooth the ground truth dot annotations can hurt the generalization bound of a model when testing on the real ground truth data. Instead, we consider crowd counting as a distribution matching problem and propose DM-Count, based on Optimal Transport, to address this problem. Unlike prior work, DM-Count does not need a Gaussian kernel to smooth the annotated dots. The generalization error bound of DM-Count is tighter than that of the Gaussian smoothed methods. Extensive experiments on four crowd counting benchmarks demonstrated that DM-Count significantly outperforms previous state-of-the-art methods.

## Broader Impact

Our work is able to more accurately estimate the crowd size in images or videos, such that it can guide crowd control and improve public safety. The estimated crowd count results are interpretable, with better crowd localization, which will increase transparency of the results for critical applications. In an age when the size of the crowd in various political events often becomes a point of heated dispute, having transparent, accurate and objective counting methods could help the historical record, as well a public acceptance of the estimates. Our method could potentially be used to protect public health by monitoring social distancing which is becoming increasingly important during the current epidemic. This method does not leverage biases in the data. The proposed method for counting is general, with possible applications to biomedical cell counting, live stock counting and etc. Our work can be adapted to count moving crowds.

## Acknowledgements

This research was partially supported by US National Science Foundation Award IIS-1763981, the SUNY2020 Infrastructure Transportation Security Center, and Air Force Research Laboratory (AFRL) DARPA FA8750-19-2-1003, the Partner University Fund, and a gift from Adobe.

## References

- [1] Deepak Babu, Neeraj Sajjan, VR Babu, and Mukundhan Srinivasan. Divide and grow: capturing huge diversity in crowd images with incrementally growing cnn. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [2] Deepak Babu Sam, Shiv Surya, and R. Venkatesh Babu. Switching convolutional neural network for crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [3] Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [4] Xinkun Cao, Zhipeng Wang, Yanyun Zhao, and Fei Su. Scale aggregation network for accurate and efficient crowd counting. In *Proceedings of the European Conference on Computer Vision*, 2018.
- [5] Antoni B Chan and Nuno Vasconcelos. Bayesian poisson regression for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2009.
- [6] Ke Chen, Chen Change Loy, Shaogang Gong, and Tony Xiang. Feature mining for localised crowd counting. In *Proceedings of the British Machine Vision Conference*, 2012.
- [7] Ke Chen, Shaogang Gong, Tao Xiang, and Chen Change Loy. Cumulative attribute space for age and crowd density estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013.
- [8] Zhi-Qi Cheng, Jun-Xiu Li, Qi Dai, Xiao Wu, and Alexander G. Hauptmann. Learning spatial awareness to improve crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [9] Pavel Dvurechensky, Alexander Gasnikov, and Alexey Kroshnin. Computational optimal transport: Complexity by accelerated gradient descent is better than by sinkhorn’s algorithm. In *Proceedings of the International Conference on Machine Learning*, 2018.
- [10] Junyu Gao, Qi Wang, and Xuelong Li. Pcc net: Perspective crowd counting via spatial convolutional network. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [11] Junyu Gao, Qi Wang, and Yuan Yuan. Scar: Spatial-/channel-wise attention regression networks for crowd counting. *Neurocomputing*, 363:1–8, 2019.
- [12] Weina Ge and Robert T Collins. Marked point processes for crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, 2014.

- [14] Haroon Idrees, Imran Saleemi, Cody Seibert, and Mubarak Shah. Multi-source multi-scale counting in extremely dense crowd images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013.
- [15] Haroon Idrees, Muhammad Tayyab, Kishan Athrey, Dong Zhang, Somaya Al-Maadeed, Nasir Rajpoot, and Mubarak Shah. Composition loss for counting, density map estimation and localization in dense crowds. In *Proceedings of the European Conference on Computer Vision*, 2018.
- [16] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [17] Xiaolong Jiang, Zehao Xiao, Baochang Zhang, Xiantong Zhen, Xianbin Cao, David Doermann, and Ling Shao. Crowd counting and density estimation by trellis encoder-decoder networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [18] Issam H Laradji, Negar Rostamzadeh, Pedro O Pinheiro, David Vazquez, and Mark Schmidt. Where are the blobs: Counting by localization with point supervision. In *Proceedings of the European Conference on Computer Vision*, 2018.
- [19] Min Li, Zhaoxiang Zhang, Kaiqi Huang, and Tieniu Tan. Estimating the number of people in crowded scenes by mid based foreground segmentation and head-shoulder detection. In *Proceedings of the International Conference on Pattern Recognition*, 2008.
- [20] Yuhong Li, Xiaofan Zhang, and Deming Chen. Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [21] Dongze Lian, Jing Li, Jia Zheng, Weixin Luo, and Shenghua Gao. Density map regression guided detection network for rgb-d crowd counting and localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [22] Sheng-Fuu Lin, Jaw-Yeh Chen, and Hung-Xin Chao. Estimation of number of people in crowded scenes using perspective transformation. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, 31(6):645–654, 2001.
- [23] Chenchen Liu, Xinyu Weng, and Yadong Mu. Recurrent attentive zooming for joint crowd counting and precise localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [24] Jiang Liu, Chenqiang Gao, Deyu Meng, and Alexander G Hauptmann. Decidenet: Counting varying density crowds through attention guided detection and density estimation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [25] Lingbo Liu, Zhilin Qiu, Guanbin Li, Shufan Liu, Wanli Ouyang, and Liang Lin. Crowd counting with deep structured scale integration network. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [26] Ning Liu, Yongchao Long, Changqing Zou, Qun Niu, Li Pan, and Hefeng Wu. Adcrowdnet: An attention-injective deformable convolutional network for crowd understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [27] Weizhe Liu, Mathieu Salzmann, and Pascal Fua. Context-aware crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [28] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. Leveraging unlabeled data for crowd counting by learning to rank. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7661–7669, 2018.
- [29] Yuting Liu, Miaoqing Shi, Qijun Zhao, and Xiaofang Wang. Point in, box out: Beyond counting persons in crowds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [30] Erika Lu, Weidi Xie, and Andrew Zisserman. Class-agnostic counting. In *Proceedings of the Asian Conference on Computer Vision*, 2018.
- [31] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. Bayesian loss for crowd count estimation with point supervision. In *Proceedings of the International Conference on Computer Vision*, 2019.

- [32] SC Martin Arjovsky and Leon Bottou. Wasserstein generative adversarial networks. In *Proceedings of the International Conference on Machine Learning*, 2017.
- [33] Daniel Onoro-Rubio and Roberto J López-Sastre. Towards perspective-free object counting with deep learning. In *Proceedings of the European Conference on Computer Vision*, 2016.
- [34] Gabriel Peyré, Marco Cuturi, et al. Computational optimal transport. *Foundations and Trends® in Machine Learning*, 11(5-6):355–607, 2019.
- [35] Viresh Ranjan, Hieu Le, and Minh Hoai. Iterative crowd counting. In *Proceedings of the European Conference on Computer Vision*, 2018.
- [36] Viresh Ranjan, Boyu Wang, Mubarak Shah, and Minh Hoai. Uncertainty estimation and sample selection for crowd counting. In *Proceedings of the Asian Conference on Computer Vision*, 2020.
- [37] Deepak Babu Sam and R Venkatesh Babu. Top-down feedback for crowd counting convolutional neural network. *Proceedings of AAAI Conference on Artificial Intelligence*, 2019.
- [38] Deepak Babu Sam, Shiv Surya, and R Venkatesh Babu. Switching convolutional neural network for crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [39] Zan Shen, Yi Xu, Bingbing Ni, Minsi Wang, Jianguo Hu, and Xiaokang Yang. Crowd counting via adversarial cross-scale consistency pursuit. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [40] Miaojing Shi, Zhaohui Yang, Chao Xu, and Qijun Chen. Revisiting perspective information for efficient crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [41] Vishwanath A Sindagi and Vishal M Patel. Cnn-based cascaded multi-task learning of high-level prior and density estimation for crowd counting. In *IEEE International Conference on Advanced Video and Signal Based Surveillance*, 2017.
- [42] Vishwanath A Sindagi and Vishal M Patel. Ha-ccn: Hierarchical attention-based crowd counting network. *IEEE Transactions on Image Processing*, 29:323–335, 2019.
- [43] Vishwanath A Sindagi and Vishal M Patel. Multi-level bottom-top and top-bottom feature fusion for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [44] Yukun Tian, Yiming Lei, Junping Zhang, and James Z Wang. Padnet: Pan-density crowd counting. *IEEE Transactions on Image Processing*, 29:2714–2727, 2019.
- [45] Alexandre B Tsybakov. *Introduction to nonparametric estimation*. Springer Science & Business Media, 2008.
- [46] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- [47] Jia Wan and Antoni Chan. Adaptive density map generation for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [48] Jia Wan, Wenhan Luo, Baoyuan Wu, Antoni B. Chan, and Wei Liu. Residual regression with semantic prior for crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [49] Chuan Wang, Hua Zhang, Liang Yang, Si Liu, and Xiaochun Cao. Deep people counting in extremely dense crowds. In *Proceedings of the ACM Multimedia Conference*, 2015.
- [50] Qi Wang, Junyu Gao, Wei Lin, and Yuan Yuan. Learning from synthetic data for crowd counting in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [51] Qi Wang, Junyu Gao, Wei Lin, and Xuelong Li. Nwpu-crowd: A large-scale benchmark for crowd counting. *arXiv preprint arXiv:2001.03360*, 2020.
- [52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- [53] Haipeng Xiong, Hao Lu, Chengxin Liu, Liang Liu, Zhiguo Cao, and Chunhua Shen. From open set to closed set: Counting objects by spatial divide-and-conquer. In *Proceedings of the International Conference on Computer Vision*, 2019.

- [54] Chenfeng Xu, Kai Qiu, Jianlong Fu, Song Bai, Yongchao Xu, and Xiang Bai. Learn to scale: Generating multipolar normalized density maps for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [55] Zhaoyi Yan, Yuchen Yuan, Wangmeng Zuo, Xiao Tan, Yezhen Wang, Shilei Wen, and Errui Ding. Perspective-guided convolution networks for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [56] Anran Zhang, Jiayi Shen, Zehao Xiao, Fan Zhu, Xiantong Zhen, Xianbin Cao, and Ling Shao. Relational attention network for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [57] Anran Zhang, Lei Yue, Jiayi Shen, Fan Zhu, Xiantong Zhen, Xianbin Cao, and Ling Shao. Attentional neural fields for crowd counting. In *Proceedings of the International Conference on Computer Vision*, 2019.
- [58] Cong Zhang, Hongsheng Li, Xiaogang Wang, and Xiaokang Yang. Cross-scene crowd counting via deep convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [59] Qi Zhang and Antoni B. Chan. Wide-area crowd counting via ground-plane density maps and multi-view fusion cnns. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [60] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. Single-image crowd counting via multi-column convolutional neural network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [61] Muming Zhao, Jian Zhang, Chongyang Zhang, and Wenjun Zhang. Leveraging heterogeneous auxiliary tasks to assist crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [62] Tao Zhao and Ramakant Nevatia. Bayesian human segmentation in crowded situations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2003.