# Semi-supervised Crowd Counting via Density Agency[*]

Hui Lin
linhuixjtu@gmail.com
School of Cyber Science and
Engineering
Xi'an Jiaotong University PRC

Zhiheng Ma
zh.ma@siat.ac.cn
Shenzhen Institute of Advanced
Technology
Chinese Academy of Science PRC

Xiaopeng Hong[†]
hongxiaopeng@ieee.org
School of Computer Science and
Technology
Harbin Institute of Technology PRC

Yaowei Wang
wangyw@pcl.ac.cn
Peng Cheng Laboratory PRC

Zhou Su
zhousu@ieee.org
Xi'an Jiaotong University PRC

## ABSTRACT

In this paper, we propose a new agency-guided semi-supervised counting approach. First, we build a learnable auxiliary structure, namely the density agency to bring the recognized foreground regional features close to corresponding density sub-classes (agents) and push away background ones. Second, we propose a density-guided contrastive learning loss to consolidate the backbone feature extractor. Third, we build a regression head by using a transformer structure to refine the foreground features further. Finally, an efficient noise depression loss is provided to minimize the negative influence of annotation noises. Extensive experiments on four challenging crowd counting datasets demonstrate that our method achieves superior performance to the state-of-the-art semi-supervised counting methods by a large margin. Code is available.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; *Semi-supervised learning settings*; *Learning from implicit feedback*.

## KEYWORDS

Crowd Counting, Semi-supervised, Contrastive Learning

## 1 INTRODUCTION

Crowd counting [13, 15, 20, 43] is to estimate the number of people in an image that could be very crowded. Since its wide applications in crowd surveillance, congestion estimation and public security, it has gained considerable attention in recent years. However, a stable and accurate crowd counter often relies on sufficient labeled images, where annotations are marked at every head centers. The annotation process can be labor-intensive and time-consuming. For instance, UCF-QNRF [7] dataset involves 1.25 million annotations, which takes 2,000 human hours in total. Therefore, semi-supervised crowd counting attempts to reduce the annotation cost while achieving good counting accuracy. It takes advantage of and extracts extra knowledge from unlabeled data, which can be obtained cheaply.
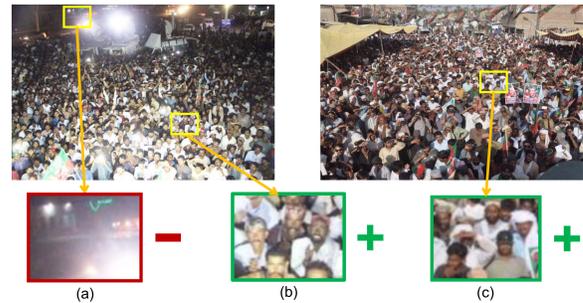


Figure 1: The density cues can be explored by correlating image regions of different count levels. Although the regions (a/b) are cropped from the same image, their semantics, i.e. the densities represented by them are quite different. Instead, we can find (c) similar to the foreground (b) in other images.

Previous semi-supervised counting methods rely on an extra auxiliary task like binary segmentation [19, 24] and count ranking [17, 18] to (self-) supervise the counting network in a multi-task paradigm. They usually generate pseudo-labels using methods like Gaussian process [32] and teacher-student models [24]. These methods can effectively reduce the annotation cost. Nonetheless, the performance is far from perfect. On the one hand, self-supervised learning leveraging rules on isolated images may ignore rich information among different images. On the other hand, it is difficult to produce high-quality pseudo-labels on top of the coarse, image-level similarity, which are calculated through neural nets trained on a small amount of labeled images.

To address the limitations, in this paper, we propose a new semi-supervised counting framework for reliable and sufficient supervision from limited labels. The novelty of this paper comes from a density agency mechanism for semi-supervised counting modeling.

The rationale of the density agency arises from the fact that the density cues can be mined by comparing image regions. An example is shown in Figure 1. When we compare a background region (a) and a foreground one (b), though the two regions are from the same image, they appear quite different. On the other hand, for a foreground region (b), we can easily find its apparent counterpart with a similar count like (c) in other crowd images. As there are huge numbers or even infinite background regions like (a), it is impossible to model them all. A more plausible alternative is to model foreground instances in line with different count levels.

On this basis, we propose the density agency guided semi supervised crowd counting model to associate the regions of different images for obtaining reliable supervised signals. The *density agency* is a learnable auxiliary structure associated with the counting backbone network. It is made up of a group of density-level agents and an agent allocator, and acts as an intermediary to identify whether an instance (an image region) is foreground or not. For foreground instances, it further quantifies the density levels of foreground instances, subdivides the foreground instances (regions with crowds) into sub-classes according to their density levels, and assigns corresponding density-level agents. It then measures the correlations between the instances and agents, which plays a significant role in guiding the learning process of the crowd counter. Generally speaking, during training, the density agency brings the recognized foreground instances to the corresponding sub-classes (density-level agents) and pushes away negative instances in feature space. More specifically, we design an efficient learning method to optimize the agents. We then propose a density-level guided contrastive loss to provide fine-grained signals for consolidating the learning of feature extractors. We further introduce a transformer structure to bridge the feature extractor and the density estimation module and refine the recognized foreground features. Finally, we devise a noise depression Bayesian loss to eliminate the annotation noises.

We evaluate the proposed model on four challenging crowd counting datasets, i.e. ShanghaiTech A and B [51], UCF-QNRF [7] and JHU-Crowd++ [30]. The experimental results show that no matter what percentage of data is labeled, our model outperforms state-of-the-art semi-supervised counting methods by a large margin. For example, under the most challenging setting of only 5% labeled data on the QNRF dataset, our method sets up new state-of-art accuracy and reduces the MAE from 160.0 to 120.2, achieving a notable about 25% reduction in mean absolute error.

In summary, we propose a novel semi-supervised crowd counting model, named Density Agency based Crowd Counting (DACount), with the following main contributions:

- We introduce a density agency guided semi-supervised learning scheme, which is able to construct sufficient supervisions for unlabeled data, breaking the supervision border between images.
- We design a density-guided uncertainty-aware contrastive loss based on the semantic differences of foreground features.
- We design a noise depression Bayesian loss to alleviate the negative influence of annotation noise.
- We set up new state-of-the-art performance for semi-supervised crowd counting on popular benchmarks.

## 2 RELATED WORKS

### 2.1 Fully-supervised Crowd Counting

In recent years, fully-supervised crowd counting has experienced a rapid development. With the help of deep CNN, the multi-column model [51] is proposed to regress the prediction to the pseudo density map generated with adaptive Gaussian kernels. CSRnet uses dilated kernels to perform accurate count estimation of highly congested scenes [11]. In [2], a scale aggregation network is proposed to exploit the local correlation and structural similarity. Ma *el al.* adopt Bayesian assumption and introduces Bayesian Loss (BL) to calculate

the expected count of pixels [21]. Wang *el al.* propose to use optimal transport to match the distributions of the point annotations and the density maps [38]. In [23], the unbalanced optimal transport is introduced to quantify the discrepancy between predicted density maps and point annotations. S3 further proposes semi-balanced Sinkhorn divergence to solve the limitations of amount constraint and entropic bias [12]. In addition, methods based on multi-scale mechanisms [22, 31, 50], segmentation [26] and perspective estimation [28, 48] are proposed to overcome the limitations of perspective and scale variations in crowd images. Moreover, BinLoss [29], P2PNet [34], UEPNet [39] and MFDC [16] further propose new solutions for crowd counting. Lately, the development of Transformer [37] also boosts counting performance. The study [45] learns adaptive feature representations with a deformable attention in Transformer network. In [15], a learnable local region attention with a corresponding regularization is proposed to address large-scale variations in crowd images.

### 2.2 Semi-/Weakly supervised Crowd Counting

As labeling crowd images is time-consuming and labor-intensive, researchers gradually focus on semi-/weakly supervised learning to find low-cost solutions for crowd counting. L2R introduces a crop ranking loss by learning containment relationships to exploit unlabeled images [17]. Meanwhile, its extension improves the ranking strategy and regards it as a self-supervised proxy task [18]. Yang *et al.* [49] propose a soft-label sorting network and only uses the crowd number as supervision. The study [32] proposes a Gaussian Process-based iterative learning mechanism to generate pseudo-ground truth for unlabeled data. Furthermore, IRAST leverages on a set of inter-related binary segmentation tasks for self-training to exploit the underlying constraints of unlabeled data [19]. The spatial uncertainty aware teacher-student framework is proposed to alleviate the noisy supervision from unlabeled data [24].

Our method is distinct to previous semi-supervised counting approaches. The methods [24, 32] generate pseudo-labels based on image-level similarity, while our method takes direct advantage of the region-level features. Moreover, compared to self-supervision methods which usually work with a single image [17–19], we make use of much richer supervision signals among different images.

### 2.3 Contrastive Learning

Contrastive learning can be formally defined by creating positive and negative data pairs, which are semantically similar and dissimilar respectively. The goal is to pull positive pairs together and push negative pairs apart in the feature space [5]. It has been widely used in self-supervised representation learning [3, 4, 6, 27], where the contrastive loss serves as an unsupervised objective function to measure the similarities of sample pairs. Previous works also apply contrastive learning into various vision tasks, such as object detection [35, 42, 44, 47] and semantic segmentation [1, 41, 52]. These tasks inherently have object/class level positive and negative samples/features, which suit for modeling the contrast. Recently, some works also adopt contrastive idea in semi-supervised learning [1, 8, 10, 33].

The objective of contrastive learning can be formulated in various functions, including max-margin contrastive loss [5], triplet
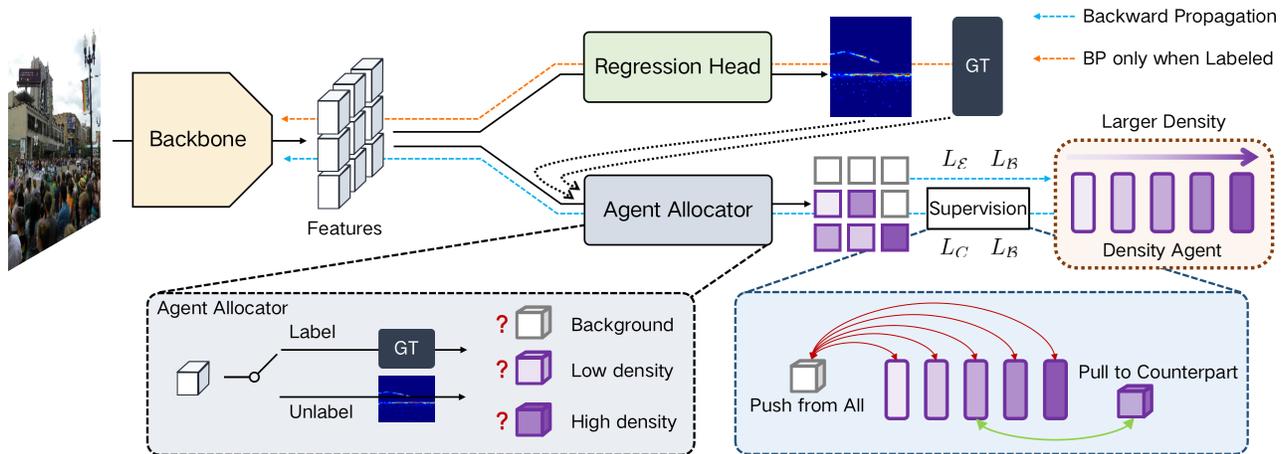
**Figure 2: Illustration of our counting framework and the training pipeline of the density agency. Features extracted by the backbone will be divided into background and different subclasses of foreground by the agent allocator. Then the agents directly supervise the backbone at intermediate feature layers instead of predicted maps by pushing all background farther and pulling the right foreground pairs closer. We design specific loss functions to update the density agents (Eq. 2 and Eq. 3) and the counting model (Eq. 7 and Eq. 3) respectively. The direction of the blue arrow represents the back-propagation update process of the loss gradients.**

loss [46], and angular loss [40]. In this work, we mainly consider the InfoNCE loss [25] and revisit contrastive learning in the context of the semi-supervised dense regression problem. We design a density guided agent and mine the relationships between positive and negative features, which yields improved performance.

## 3 THE PROPOSED METHOD

DACount is a combination of four novel modules to semi-supervised crowd counting: a learnable density agency, a self-supervised contrastive loss, a foreground transformer and a noise depression loss. In this section, we first introduce the notations for the semi-supervised counting problem and then introduce a density agency to establish a shared supervision scheme among all training images. It consists of an agent allocator and a set of density agents. On that basis, we propose the density agent guided contrastive learning loss. Then we elaborate the proposed network structure with a transformer to rectify foreground features. Finally, a noise depression Bayesian loss is introduced to enhance the update of regression head under few labeled data.

### 3.1 Density Agent

Assume we have a labeled dataset $\mathcal{L} = \{(x^i, y^i)\}_{i=1}^{N_l}$ of $N_l$ labeled samples, where $x^i$ is a crowd image and $y^i$ is the corresponding point ground truth. An unlabeled dataset can also be denoted as $\mathcal{U} = \{(x^i)\}_{i=1}^{N_u}$ of $N_u$ unlabeled samples, where the size of it is usually much larger than that of the $\mathcal{L}$, i.e., $N_l << N_u$.

In crowd counting, as shown by the example in Figure 1, we discover that the region-level variations among different images are mainly from background, while regions with crowds may have similar appearance. Therefore, based on this observation, we propose to build a density agency targeted on foreground across all images. The agency is a learnable auxiliary structure associated with the counting model, acting as an intermediary to explore the

correlations of foreground from different images and background in feature space. Specifically, the agent consists of a group of density-level agents and an agent allocator.

The **agent allocator** is designed to divide features extracted by the backbone into background and different groups of foreground features in line with the quantitative levels of their density and assign **density agents** accordingly. It is formulated as follows. Suppose the density range is partitioned into $N_a$ intervals, the borders of which are denoted by a series, i.e., $(0, v_1, v_2, \ldots, v_{N_a-1})$. We have $N_a$ density-level agents $\mathcal{A} = \{f_i\}_{i=1}^{N_a}$, $f_i \in \mathbb{R}^c$, where $c$ is the channel dimension. The agents can be assigned by the following semantic-driven rule.

$$f_1 \leftarrow (0, v_1), \; f_2 \leftarrow [v_1, v_2), \; \ldots, \; f_{N_a} \leftarrow [v_{N_a-1}, +\infty).$$

The left-bottom rectangle in Figure 2 illustrates how the agent allocator maps the regions from an input image to the corresponding agents. For labeled data, the agent allocator will leverage on the ground-truth to quantify the density levels and classify different regions. For unlabeled data, there is no ground truth available. To obtain the region-wise density levels, the unlabeled image first go through the backbone and the regression head, and then the agent allocator use the predicted density map as the reference.

More specifically, let $\mathcal{E}$ and $\mathcal{B}$ denote the foreground regional feature set and the background regional feature set respectively. $\mathcal{E} \cap \mathcal{B} = \varnothing$, $\mathcal{E} \cup \mathcal{B} = \mathcal{F}$, where $\mathcal{F}$ denotes all region features extracted by the backbone. For a foreground feature $e \in \mathcal{E}$, if the density value of $e$ is in the range $[v_i, v_{i+1})$, we find its agent $f_i = Z(e)$ for that interval through the agent allocator.

With the agency mechanism, we expect to pull the foreground regional features close to their density agents for the corresponding density intervals and push the background features away from all the agents. As a result, an agent is desired to uniquely represent a foreground density interval and be significantly different to the regional background features. To achieve this goal, we design a

simple yet efficient agent learning algorithm, which is shown in the right-bottom part of Figure 2.

Given $\mathcal{E}$ and $\mathcal{B}$, we optimize the agents to maximize the similarity between positive pairs, each of which is formed of a foreground feature vector and its density agent and minimize the one between negative pairs, each of which contains a background feature vector and any agents. In details, the process is formulated as follows.

$$\min\left(L_{\mathcal{E}} + L_{\mathcal{B}}\right), \tag{1}$$

where

$$L_{\mathcal{E}} = -\sum_{e \in \mathcal{E}} s(Z(e), e) \tag{2}$$

and

$$L_{\mathcal{B}} = \frac{1}{N_a} \sum_{f \in \mathcal{A}} \sum_{b \in \mathcal{B}} s(f, b). \tag{3}$$

$L_{\mathcal{E}}$ is responsible for pulling close positive pairs, i.e., $(Z(e), e)$ and $L_{\mathcal{B}}$ for pushing away any agent $f$ and a background feature $b$.

By taking into account the range of values, we adopt the cosine similarity between two vectors as the measurement.

$$s(f_1, f_2) = \frac{f_1 \cdot f_2}{\|f_1\|\|f_2\|}. \tag{4}$$

At every stage, we use gradient descent to update the agents in order to reach the optimal solution for the current state. The optimization process of foreground and background loss can be written by,

$$
\begin{aligned}
f_{t+1} &= f_t - \gamma(\nabla_f L_{\mathcal{E}} + \nabla_f L_{\mathcal{B}}), \text{ where} \\
\nabla_f L_{\mathcal{E}} &= \sum_{Z(e)=f} s(f, e)\frac{f}{\|f\|^2} - \frac{e}{\|e\|\|f\|}, \\
\nabla_f L_{\mathcal{B}} &= \sum_{b \in \mathcal{B}} \frac{b}{\|b\|\|f\|} - s(f, b)\frac{f}{\|f\|^2}.
\end{aligned}
\tag{5}
$$

$\gamma$ is the learning rate. For computational simplicity and efficiency, the Adam optimizer [9] is used to update the density agents.

## 3.2 Density Agent Guided Contrastive Learning

Given a foreground feature, either labeled or unlabeled, the agent allocator assigns an agent to it. However, there are two key problems. First, as there are errors in the generated density maps (for unlabeled data), the allocator inevitably encounters uncertainties in matching agents and features. As a result, a wrong pulling will further lead to poorer density prediction. Second, Eq. 2 and 3 only supervise the differences between foreground and background, and does not make a finer distinction among foreground densities. As it shown in Figure 3 (b), this may make the agent collapse to a common position in the feature space. In this section, our goal is to build a reliable scheme to first select signals from foreground supervision and then leverage on the contrastive learning to separate different foreground features.

In practice, especially in noisy cases when only partial data are labeled, the predicted density value $d$ of a foreground feature may perturb within a certain range. Put another way, this uncertainty stands for the matching probability of a foreground feature over different agents as well. The higher the matching probability, the more likely the foreground features can be regarded as a positive proposal, and vice versa. Based on this understanding, we assume
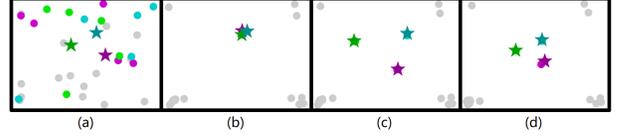


(a)  (b)  (c)  (d)

Figure 3: A toy experiment to show the influence of different supervisions. Colored and grey dots represent different foreground features and background respectively while stars represent the agents. (a) Initialization. (b) When using $L_{\mathcal{E}}$ in Eq. 2 instead of contrastive loss $L_c$ in Eq. 7 to update the counting model, the foreground features collapse to be similar. (c) When using $L_c$ instead of $L_{\mathcal{E}}$ to update the agency, the distance between the foreground features may be greater than that between foreground and background. (d) When using $L_{\mathcal{B}}$ in Eq. 3 and $L_c$ to update the counting model and using $L_{\mathcal{B}}$ and $L_{\mathcal{E}}$ to update the agency, foreground features are similar enough and also clearly differentiated.

that this uncertainty conforms to a Laplace distribution and assign higher weights to proposals that are clearly divided, and down weigh those which are on the boundary between positive and negative. As a result, suppose $d$ falls into the $i$-th interval, we define the uncertainty-aware weight $\omega_i$ as,

$$\omega_i = 8|\hat{\omega}_i - 0.25|, \tag{6}$$

where $\hat{\omega}_i$ is the matching probability, subject to a Laplace distribution, i.e., $\hat{\omega}_i \sim \textbf{Laplace}\,(a_i, 2) = \frac{1}{4}e^{-\frac{|d-a_i|}{2}}$. The location parameter $a_i$ is the central value of the density interval for $f_i$. In extreme cases, $a_1 = \frac{1}{2}v_1$ and $a_{N_a} = v_{N_a-1}$. The parameters in Eq. 6 are set to ensure the value of the weight $\omega_i$ within the range of $[0, 1]$.

In light of Eq. 6, we select reliable pairs for supervision. We consider features which the matching probability above the certain threshold 0.25 to be positive proposals and others to be negative. Matching positive pairs should be pulled in further, while selected negatives should be pushed farther, aiming to increase the dissimilarity of features with different density values. Then for feature $e$, we define the uncertainty-aware contrastive loss as

$$L_c(e) = -\log \frac{\sum_{\hat{\omega}_i \geq 0.25} \omega_i u_i}{\sum_{j=1}^{N_a} \omega_j u_j}, \quad u_i = e^{s(f_i, e)/\tau} \tag{7}$$

where $\tau$ is a temperature parameter. Then the final loss is the sum of contrastive loss of all foreground features and background loss with weight $\lambda_b$:

$$L_C = \sum_{e \in \mathcal{E}} L_c(e) + \lambda_b L_{\mathcal{B}}. \tag{8}$$

In experiments, we will set $\lambda_b = 1$ for simplification.

Specifically, we use different loss functions to update the density agency and the counting model respectively. To update the density agents, we adopt Eq. 2 as foreground loss and Eq. 3 as background loss. The detailed learning process is described in the last few paragraphs of Section 3.1. Meanwhile, to update the counting model, we adopt the contrastive loss in Eq. 7 as foreground loss and Eq. 3 as background loss. This update solution is in consideration of the pace of gradient updates. For a negative counterpart $f_i$, $\hat{\omega}_i < 0.25$,
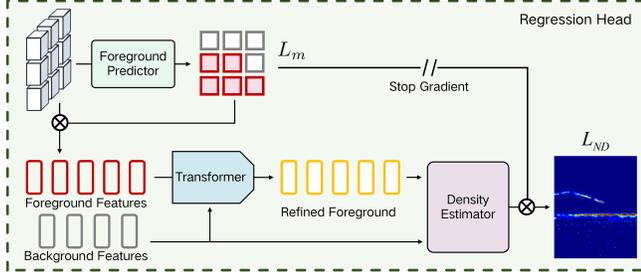
Figure 4: Details of the regression head. It consists of a foreground predictor to segment foreground and background, a transformer to refine foreground features and a density estimator to generate the density map. It is trained only using labeled data.

its gradient of Eq. 7 is

$$\nabla_f L_c = \frac{1}{\tau} \cdot \frac{\omega_i u_i}{\sum_{j=1}^{N_a} \omega_j u_j} \cdot \left( \frac{e}{\|e\| \|f_i\|} - s(f_i, e) \frac{f_i}{\|f_i\|^2} \right). \quad (9)$$

Comparing to the gradient of background loss in Eq. 5, an inappropriate negative pair with a larger similarity may cause the agent to be pushed away by a pace greater than that of the background. As shown in Figure 3 (c), the distance between the foreground may be greater than that between the foreground and background features when using the contrastive loss for agent update.

## 3.3 Network Structure

To support the above functions, a powerful and reliable regression head is essential. In this regard, we build the semi-supervised training scheme and further introduce a transformer for refining the foreground features.

For the training scheme, we decide to only leverage on labeled data to train the regression head in order to avoid contamination by inaccurate supervision from unlabeled data, keeping a 'pure' head. To achieve this, the unlabeled data will only act directly on the feature extractor, instead of the whole model. Therefore, when using the predicted density map to divide features, we will stop gradients of the regression head.

The total framework of our counting model is shown in Figure 4. For an input image, we first use a CNN as the backbone to extract features. Then the regression head consists of three modules: a foreground predictor, a transformer and a density estimator, all of which are only trained by labeled data to remain 'pure'.

The extracted features are first transmitted into the foreground predictor to get a division mask $M \in \{0, 1\}^N$ where $N$ is the number of pixels. Then the foreground and background features can be divided by

$$\mathcal{E} \leftarrow \text{sg}(M) \circ \mathcal{F}, \quad \mathcal{B} \leftarrow \text{sg}(1 - M) \circ \mathcal{F}. \quad (10)$$

sg denotes stop gradients and $\circ$ denotes Hadamard product. The predicted foreground and background mask will be supervised by the generated *binary* foreground mask $\hat{M}$ under MSE loss, which is defined by

$$L_m = \sqrt{\sum (\hat{M} - M)^2}. \quad (11)$$

Next, we use the transformer structure to refine only the foreground features, which detail is shown in Figure 5. It includes a
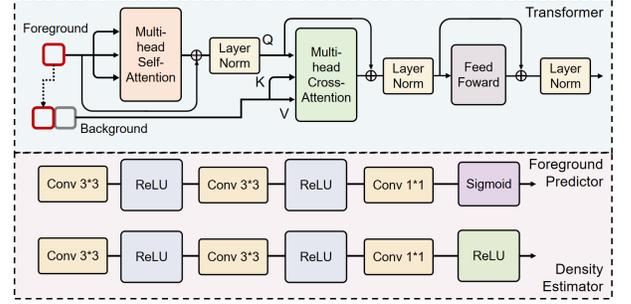


Figure 5: Detailed structures of the transformer, foreground predictor and density estimator.

self-attention module for foreground and a cross-attention between foreground and background features. With the help of these two modules, the foreground features can not only generate interactive attention with themselves, but also mine the relations with the background. The distance between the foreground and background features is further enlarged, and the differences within the foreground are also fully explored.

Finally, the background features and the refined foreground features are fed together into the density estimator to predict the final density map.

## 3.4 Noise Depression Bayesian Loss

To ensure the quality of the regression head, only labeled data are used to train it. However, there are still location noises in the labels. Here, we propose a noise depression loss to minimize the negative influence to regression head by these noises.

We believe that when the gap between ground truth value and expected value the model predicted is too large, there may be noise in the label. Therefore, we propose to add a modulating factor to down-weigh these large gaps. Given the predicted density map D, the noise depression Bayesian loss is defined by

$$\varepsilon_j = |\delta(y_j) - \sum_i^N p_{ij} \cdot D_i|,$$

$$L_{ND} = \sum_{j=1}^M e^{-\beta \text{sg}(\varepsilon_j)} \cdot \varepsilon_j. \quad (12)$$

sg($\cdot$) means stop gradients and $p_{ij}$ represents the contributed posterior probability of ground-truth $j$ for predicted pixel $i$. $\beta \geq 0$ is the depression parameter. When $\beta = 0$, the loss will degenerate into normal Bayesian loss. Increasing the value of $\beta$, the loss will gradually ignore the penalty for large gap. we found $\beta = 1$ to work best in our experiments. Then the total loss under labeled data turns into

$$L_{label} = L_{ND} + \lambda_m L_m + \lambda_c L_C, \quad (13)$$

while the loss under unlabeled data is

$$L_{unlabel} = \lambda_c L_C. \quad (14)$$

And the total loss of the whole model is a combination of $L_{label}$ and $L_{unlabel}$ with an unlabeled parameter $\lambda_u$ is

$$L = L_{label} + \lambda_u L_{unlabel}. \quad (15)$$

The training process of DACount is summarised in Algorithm 1.

**Algorithm 1:** DACount Learning

---

**Input:** Sampled labeled images $\mathcal{X}_l \in \mathcal{L}$ and unlabeled
    images $\mathcal{X}_u \in \mathcal{U}$.

1 **Output:** The counting model $R$ which consists of feature
    backbone $g$ and regression head $h$, $R = h \circ g$. The
    density agent $\mathcal{A}$.

2 Initialize $L_{\mathcal{A}} \leftarrow 0, L \leftarrow 0$;

3 **for** $x \in \mathcal{X}_l \cup \mathcal{X}_u$ **do**

4    Get features $\mathcal{F} \leftarrow g(x)$;

5    Get foreground mask and predicted density map
      $M, D \leftarrow h(\mathcal{F})$;

6    **if** $x \in \mathcal{X}_l$ **then**

7       Generate ground-truth foreground mask $\hat{M}$ and
         ground-truth density map $\hat{D}$;

8       $L \leftarrow L + L_{ND}$ based on Eq. 12;

9       $L \leftarrow L + \lambda_m L_m$ based on MSE loss between $M$ and
         $\hat{M}$;

10      Update regression head $h$ minimizing $L$;

11   **end**

12   Get foreground and background features based on
      Eq. 10 by $M$ or $\hat{M}$;

13   Divide $\mathcal{E}$ by $D$ or $\hat{D}$;

14   Calculate $L_{\mathcal{E}}$ and $L_{\mathcal{B}}$ based on Eq. 2 and Eq. 3;

15   Calculate $L_C$ based on Eq. 7 and Eq. 8;

16   $L \leftarrow L + \lambda_c L_C$;

17   $L_{\mathcal{A}} \leftarrow L_{\mathcal{A}} + L_{\mathcal{E}} + L_{\mathcal{B}}$;

18   **if** $x \in \mathcal{X}_u$ **then**

19      Adjust the loss weight by $L \leftarrow \lambda_u L, L_{\mathcal{A}} \leftarrow \lambda_u L_{\mathcal{A}}$;

20   **end**

21 **end**

22 Update $\mathcal{A}$ minimizing $L_{\mathcal{A}}$;

23 Update feature backbone $g$ minimizing $L$;

24 **Return** the model $R = h \circ g$, the density agent $\mathcal{A}$.

---

# 4 EXPERIMENTS

## 4.1 Implementation Details:

**Network Details:** VGG-19 is pre-trained on ImageNet and adopted
as our CNN backbone. The number of density agents is set as
$N_a = 24$. We use Adam algorithm [9] to optimize the model and
the agent. Both set the learning rate as $10^{-6}$, specifically, $\gamma = 10^{-6}$.
For the loss parameters, we set $\lambda_m = 0.1$, $\lambda_c = 0.01$, and $\lambda_u = 0.1$.
The depression parameter $\beta = 1$. The temperature $\tau = 0.1$.

**Training Details:** We adopt random scaling of [0.7, 1.3] and hori-
zontal flipping for each training image, and use random crop with
a size of $512 \times 512$. The only exception is to use the size of $256 \times
256$ for extremely small images on ShanghaiTech A. We limit the
shorter side of each image within 2048 pixels in all datasets.

## 4.2 Datasets

We conduct extensive experiments on four crowd counting bench-
marks to verify the effectiveness of our proposed method.

**ShanghaiTech A** [51] contains 482 crowd images with 244,167
annotated points. The training set has 300 images, and the testing
set has the remaining 182 images.

**ShanghaiTech B** [51] contains 716 crowd images, which are taken
in the crowded street of Shanghai. The training set has 316 images,
and the testing set has the remaining 400 images.

**UCF-QNRF** [7] includes 1535 high-resolution images with 1.25
million annotated points. There are 1,201 images in the training set
and 334 images in the testing set.

**JHU-Crowd++** [30] includes 4,372 images with 1.51 million anno-
tated points. There are 2272 images used for training, 500 images
for validation, and the rest 1600 images used for testing.

## 4.3 Experimental Results

We evaluate DACount and compare it with state-of-the-art semi-
supervised methods, including mean-teacher (MT) [36], Learn-
ing to Rank (L2R) [17], Gaussian Process (GP) [32], IRAST [19],
IRAST+SPN (Scale Pyramid Network) [19], and SUA [24].

The results are shown in Table 1. It can been easily observed
that DACount outperforms other state-of-the-art semi-supervised
methods by a significant accuracy improvement on all datasets,
regardless of the settings of the ratio of labeled data. Under the
(relatively) easiest setting (i.e., the setting of 40% labeled data),
our method reduces the MAE by 39.2, 15.6, 1.0, and 4.5 points on
databases QNRF, JHU++, SHanghaiTech A and B respectively, com-
pared to the second best method SUA. The improvement becomes
even more apparent when there is less labeled data available. For
example, on UCF-QNRF dataset, DACount reduces all the previous
methods by at least 39.8 and 62.6 in terms of MAE and MSE, respec-
tively. These excellent results demonstrate the effectiveness of our
method in semi-supervised crowd counting.

## 4.4 Ablation study

We conduct ablation experiments to verify the effectiveness of the
three component in DACount, namely the learnable density agent,
contrastive loss, and noise depression loss, as shown in Table 2.

All experiments are conducted on UCF-QNRF with a labeled
ratio of 5%. We start with the baseline of BL [21], which only uses
labeled data and vanilla CNN to extract features. Then, we adopt
transformer module with foreground and background separated,
which detail can be referred to Section 3.3, and the performance
improves by a large margin even still without unlabeled data. We
believe that this improvement is due to the global attention of fore-
ground on the one hand, and the cross attention with background
on the other hand.

Next, we adopt the proposed density agency to study the help
of this supervision on unlabeled data. With the agency, we can
take full advantage of the potential of unlabeled data. The learnable
structure builds a bridge between labeled and unlabeled data, ex-
ploiting the correlations among different images. The error reduces
12.9 and 25.3 for MAE and MSE, respectively. Then by transferring
density information to features and widening the differences be-
tween foreground, the contrastive learning further improves the
counting accuracy. In addition, we also take into account the un-
certainty of the predicted density map, weighting the positive and
negative counterparts, and have a pleasing performance. Finally,

| Methods | Labeled Percentage | UCF-QNRF | | JHU++ | | ShanghaiTech A | | ShanghaiTech B | |
|---|---|---|---|---|---|---|---|---|---|
| | | MAE | MSE | MAE | MSE | MAE | MSE | MAE | MSE |
| MT [36] | 5% | 172.4 | 284.9 | 101.5 | 363.5 | 104.7 | 156.9 | 19.3 | 33.2 |
| L2R [17] | 5% | 160.1 | 272.3 | 101.4 | 338.8 | 103.0 | 155.4 | 20.3 | 27.6 |
| GP [32] | 5% | 160.0 | 275.0 | - | - | 102.0 | 172.0 | 15.7 | 27.9 |
| DACount (Ours) | 5% | **120.2** | **209.3** | **82.2** | **294.9** | **85.4** | **134.5** | **12.6** | **22.8** |
| MT [36] | 10% | 145.5 | 250.3 | 90.2 | 319.3 | 94.5 | 156.1 | 15.6 | 24.5 |
| L2R [17] | 10% | 148.9 | 249.8 | 87.5 | 315.3 | 90.3 | 153.5 | 15.6 | 24.4 |
| IRAST [19] | 10% | - | - | - | - | 86.9 | 148.9 | 14.7 | 22.9 |
| IRAST+SPN [19] | 10% | - | - | - | - | 83.9 | 140.1 | - | - |
| DACount (Ours) | 10% | **109.0** | **187.2** | **75.9** | **282.3** | **74.9** | **115.5** | **11.1** | **19.1** |
| MT [36] | 40% | 147.2 | 249.6 | 121.5 | 388.9 | 88.2 | 151.1 | 15.9 | 25.7 |
| L2R [17] | 40% | 145.1 | 256.1 | 123.6 | 376.1 | 86.5 | 148.2 | 16.8 | 25.1 |
| GP [32] | 40% | 136.0 | - | - | - | 89.0 | - | - | - |
| IRAST [19] | 40% | 138.9 | - | - | - | - | - | - | - |
| SUA [24] | 40% | 130.3 | 226.3 | 80.7 | 290.8 | 68.5 | 121.9 | 14.1 | 20.6 |
| DACount (Ours) | 40% | **91.1** | **153.4** | **65.1** | **260.0** | **67.5** | **110.7** | **9.6** | **14.6** |

Table 1: Comparisons with the state of the arts semi-supervised counting methods on ShanghaiTech A, ShanghaiTech B, UCF-QNRF, and JHU-Crowd++. The best performance is shown in bold. The results of other methods under the 40% labeled setting are referred to [24] and all other results are from the original papers.



GT: 137    GT: 337    GT: 349    GT: 558

Labeled Only: 2963    Labeled Only: 617    Labeled Only: 131    Labeled Only: 345

DACount: 174    DACount: 323    DACount: 345    DACount: 534
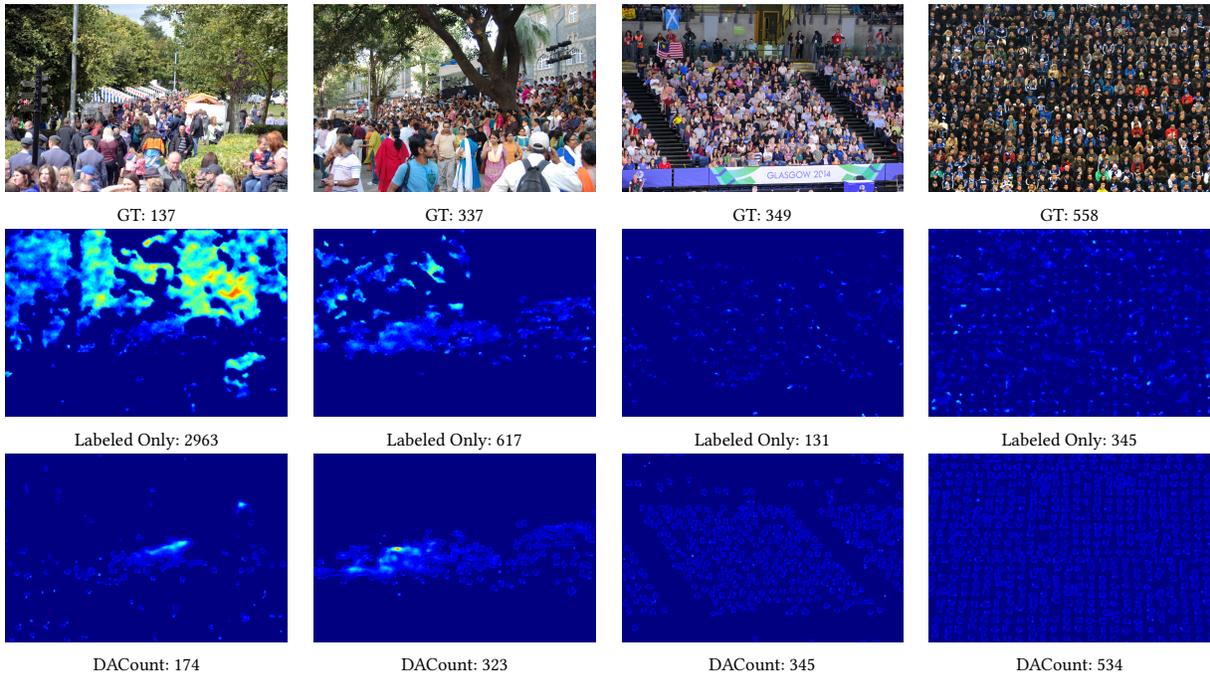
Figure 6: Visualizations of predicted densities on unlabeled training images. The first row: input images. The second row: predicted density maps trained only by labeled data. The third row: predicted density maps trained by both labeled and unlabeled data. Both models are trained on UCF-QNRF with a labeled ratio of 5%. When the model is trained only with labeled data, serious false alarms in the background (the first two examples) and apparent underestimation of the foreground density (the last two examples) are observed in the second row. In contrast, our agency-guided semi-supervised counting scheme can well exploit rich supervised signal from unlabeled data and thus produce density maps of much higher quality in the third row.

the noise depression Bayesian loss down-weights the supervision of noisy instances. The best performance is achieved when combining all above components, improving 31.2 and 54.6 for MAE and MSE compared to the baseline.

## 4.5 Discussions

**The impact of unlabeled data.** We conduct experiments to examine the impact of unlabeled data. Specifically, excluding unlabeled data can be equal to only using the labeled loss $L_{label}$. And the

| Components | MAE | MSE |
|---|---|---|
| Baseline (BL) | 151.4 | 263.9 |
| + Transformer | 141.3 | 241.2 |
| + Momentum Density Agent | 130.6 | 218.3 |
| + Learnable Density Agent | 128.4 | 215.9 |
| + Contrastive Learning | 124.7 | 214.6 |
| + Uncertainty Aware Contrastive | 122.0 | 212.9 |
| + Noise Depression Loss | **120.2** | **209.3** |

**Table 2: Ablation study on UCF-QNRF (labeled ratio 5%).**

| Loss | Labeled Percentage | MAE | MSE |
|---|---|---|---|
| $L_{label}$ | 5% | 138.3 | 241.1 |
| $L_{label} + \lambda_u L_{unlabel}$ | 5% | 120.2 | 209.3 |
| $L_{label}$ | 10% | 129.4 | 229.7 |
| $L_{label} + \lambda_u L_{unlabel}$ | 10% | 109.0 | 187.2 |
| $L_{label}$ | 40% | 98.4 | 182.7 |
| $L_{label} + \lambda_u L_{unlabel}$ | 40% | 91.1 | 153.4 |

**Table 3: The impact of using unlabeled data on UCF-QNRF.**

| $\lambda_u$ | 0 | 0.001 | 0.005 | 0.01 | 0.05 |
|---|---|---|---|---|---|
| MAE | 94.1 | 92.5 | 90.4 | 91.2 | 86.5 |
| MSE | 142.9 | 140.6 | 141.3 | 142.0 | 139.3 |

| $\lambda_u$ | 0.1 | 0.5 | 1 | GP [32] | |
|---|---|---|---|---|---|
| MAE | **85.4** | 87.9 | 94.2 | 102.0 | |
| MSE | 134.5 | **134.3** | 145.8 | 172.0 | |

**Table 4: The influence of unlabeled parameter $\lambda_u$ on Shang-haiTech A (labeled ratio 5%).**

combination of $L_{label}$ and $L_{unlabel}$ forms the proposed DACount. The comparison result is shown in Table 3.

Compared to using labeled data only, DACount which leverages on unlabeled data outperforms consistently with all labeled percentages. Specifically, it keeps a over 7.4% improvement with the help of unlabeled data. As the results shown, when the labeled percentage is smaller, the improvement will be larger. This is reasonable since the model can utilize more unlabeled data. This quantitative experiment proves that our model can make full use of unlabeled data to help improve the performance.

**The influence of semi parameter $\lambda_u$.** $\lambda_u$ controls the supervised proportion of unlabeled data. We compare the counting accuracy for using different $\lambda_u$ during training. The result is shown in Table 4.

All experiments are conducted on ShanghaiTech A with a labeled ratio of 5% and we choose previous state-of-the-art GP [32] for comparison. In general, when tuning the unlabeled parameter, the final performance has a small fluctuation and is consistently better than previous method. In specific, when $\lambda_u = 1$ and $N_u = 19N_l$, due to the lack of explicit label supervision, the accuracy is not satisfactory. And when $\lambda_u \leq 0.01$, the unlabeled loss is too small to fully help the supervision, so the performance also drops, which is almost on par with only using labeled data. During other experiments, we will choose $\lambda_u = 0.1$.

**The impact of distribution assumptions of matching probability** is studied in Table 5. In addition to the Laplace distribution, we also assume a Normal distribution where the matching probability subjects to $\hat{\omega}_i \sim \textbf{Normal}\,(a_i, 0.5) = \frac{\sqrt{2}}{\sqrt{\pi}}\,e^{-2(d-a_i)^2}$. Different

| Datasets | Laplace | | Normal | |
|---|---|---|---|---|
| | MAE | MSE | MAE | MSE |
| UCF-QNRF | 120.2 | 209.3 | 117.3 | 201.1 |
| JHU++ | 82.2 | 294.9 | 83.2 | 303.7 |
| ShanghaiTech A | 85.4 | 134.5 | 85.9 | 135.3 |
| ShanghaiTech B | 12.6 | 22.8 | 12.8 | 23.3 |

**Table 5: The comparison between two distribution assumptions in contrastive learning (labeled ratio 5%).**

assumptions of matching probability have little effect on the final result and both can achieve good accuracy.

As the results shown, adopting different distribution assumptions will have little effect on the final results. Both Laplace and Gaussian assumptions can achieve good accuracy and surpass previous state-of-the-art methods by a large margin on semi-supervised crowd counting. This suggests that the model is insensitive to the setting of the contrastive weights $\omega$.

## 5 CONCLUSION

This paper proposes a novel approach for semi-supervised crowd counting and demonstrates the benefits of building regional supervision signals across images. It is achieved by a learnable density agency, which semantically connects the foreground features to the agents and rectifies the backbone feature extraction networks. A noise depression Bayesian loss is also introduced to mitigate the problem of annotation noises. The proposed agency mechanism provides rich supervised signals while the training pipeline with gradient-stop ensures that only the most reliable supervision signals are propagated. DACount achieves excellent performance upon four challenging crowd counting datasets. Extensive experiments also demonstrate its potential in reducing annotations efforts. In future, we will apply the proposed agency based semi-supervised learning scheme to other semi-supervised learning tasks.

## REFERENCES
[1] Inigo Alonso, Alberto Sabater, David Ferstl, Luis Montesano, and Ana C Murillo. 2021. Semi-supervised semantic segmentation with pixel-level contrastive learning from a class-wise memory bank. In *ICCV*.
[2] Xinkun Cao, Zhipeng Wang, Yanyun Zhao, and Fei Su. 2018. Scale aggregation network for accurate and efficient crowd counting. In *ECCV*.
[3] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR.
[4] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. 2020. Improved baselines with momentum contrastive learning. *arXiv preprint* (2020).
[5] Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *CVPR*.
[6] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *CVPR*.
[7] Haroon Idrees, Muhmmad Tayyab, Kishan Athrey, Dong Zhang, Somaya Al-Maadeed, Nasir Rajpoot, and Mubarak Shah. 2018. Composition loss for counting, density map estimation and localization in dense crowds. In *ECCV*.

[8] Li Jiang, Shaoshuai Shi, Zhuotao Tian, Xin Lai, Shu Liu, Chi-Wing Fu, and Jiaya Jia. 2021. Guided Point Contrastive Learning for Semi-supervised Point Cloud Semantic Segmentation. In *ICCV*.

[9] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint* (2014).

[10] Junnan Li, Caiming Xiong, and Steven CH Hoi. 2021. Comatch: Semi-supervised learning with contrastive graph regularization. In *ICCV*.

[11] Yuhong Li, Xiaofan Zhang, and Deming Chen. 2018. Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes. In *CVPR*.

[12] Hui Lin, Xiaopeng Hong, Zhiheng Ma, Xing Wei, Yunfeng Qiu, Yaowei Wang, and Yihong Gong. 2021. Direct measure matching for crowd counting. *IJCAI* (2021).

[13] Hui Lin, Xiaopeng Hong, and Yabin Wang. 2021. Object Counting: You Only Need to Look at One. *arXiv preprint* (2021).

[14] Hui Lin, Zhiheng Ma, Xiaopeng Hong, Yaowei Wang, and Zhou Su. 2020. Semi-supervised Crowd Counting via Density Agency. In *the Proceedings of the 30th ACM International Conference on Multimedia (MM '22)*.

[15] Hui Lin, Zhiheng Ma, Rongrong Ji, Yaowei Wang, and Xiaopeng Hong. 2022. Boosting Crowd Counting via Multifaceted Attention. In *CVPR*.

[16] Xinyan Liu, Guorong Li, Zhenjun Han, Weigang Zhang, Yifan Yang, Qingming Huang, and Nicu Sebe. 2021. Exploiting sample correlation for crowd counting with multi-expert network. In *ICCV*.

[17] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. 2018. Leveraging unlabeled data for crowd counting by learning to rank. In *CVPR*.

[18] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. 2019. Exploiting unlabeled data in cnns by self-supervised learning to rank. *IEEE TPAMI* (2019).

[19] Yan Liu, Lingqiao Liu, Peng Wang, Pingping Zhang, and Yinjie Lei. 2020. Semi-supervised crowd counting via self-training on surrogate tasks. In *ECCV*.

[20] Zhiheng Ma, Xiaopeng Hong, Xing Wei, Yunfeng Qiu, and Yihong Gong. 2021. Towards a universal model for cross-dataset crowd counting. In *ICCV*.

[21] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. 2019. Bayesian loss for crowd count estimation with point supervision. In *ICCV*.

[22] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. 2020. Learning scales from points: A scale-aware probabilistic model for crowd counting. In *ACM Multimedia*.

[23] Zhiheng Ma, Xing Wei, Xiaopeng Hong, Hui Lin, Yunfeng Qiu, and Yihong Gong. 2021. Learning to count via unbalanced optimal transport. In *AAAI*.

[24] Yanda Meng, Hongrun Zhang, Yitian Zhao, Xiaoyun Yang, Xuesheng Qian, Xiaowei Huang, and Yalin Zheng. 2021. Spatial Uncertainty-Aware Semi-Supervised Crowd Counting. In *ICCV*.

[25] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. *arXiv preprint* (2018).

[26] Muhamad Sajid, Ali Hassan, and Shoab A Khan. 2016. Crowd counting using adaptive segmentation in a congregation. In *ICSIP*.

[27] Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. 2018. Time-contrastive networks: Self-supervised learning from video. In *ICRA*.

[28] Miaojing Shi, Zhaohui Yang, Chao Xu, and Qijun Chen. 2019. Revisiting perspective information for efficient crowd counting. In *CVPR*.

[29] Sravya Vardhani Shivapuja, Mansi Pradeep Khamkar, Divij Bajaj, Ganesh Ramakrishnan, and Ravi Kiran Sarvadevabhatla. 2021. Wisdom of (Binned) Crowds: A Bayesian Stratification Paradigm for Crowd Counting. In *ACM Multimedia*.

[30] Vishwanath Sindagi, Rajeev Yasarla, and Vishal MM Patel. 2020. Jhu-crowd++: Large-scale crowd counting dataset and a benchmark method. *PAMI* (2020).

[31] Vishwanath A. Sindagi and Vishal M. Patel. 2019. Multi-level bottom-top and top-bottom feature fusion for crowd counting. In *ICCV*.

[32] Vishwanath A Sindagi, Rajeev Yasarla, Deepak Sam Babu, R Venkatesh Babu, and Vishal M Patel. 2020. Learning to count in the crowd from limited labeled data. In *ECCV*.

[33] Ankit Singh, Omprakash Chakraborty, Ashutosh Varshney, Rameswar Panda, Rogerio Feris, Kate Saenko, and Abir Das. 2021. Semi-supervised action recognition with temporal contrastive learning. In *CVPR*.

[34] Qingyu Song, Changan Wang, Zhengkai Jiang, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yang Wu. 2021. Rethinking counting and localization in crowds: A purely point-based framework. In *ICCV*.

[35] Bo Sun, Banghuai Li, Shengcai Cai, Ye Yuan, and Chi Zhang. 2021. Fsce: Few-shot object detection via contrastive proposal encoding. In *CVPR*.

[36] Antti Tarvainen and Harri Valpola. 2017. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *NIPS* (2017).

[37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.

[38] Boyu Wang, Huidong Liu, Dimitris Samaras, and Minh Hoai Nguyen. 2020. Distribution matching for crowd counting. *NIPS* (2020).

[39] Changan Wang, Qingyu Song, Boshen Zhang, Yabiao Wang, Ying Tai, Xuyi Hu, Chengjie Wang, Jilin Li, Jiayi Ma, and Yang Wu. 2021. Uniformity in Heterogeneity: Diving Deep into Count Interval Partition for Crowd Counting. In *ICCV*.

[40] Jian Wang, Feng Zhou, Shilei Wen, Xiao Liu, and Yuanqing Lin. 2017. Deep metric learning with angular loss. In *ICCV*.

[41] Wenguan Wang, Tianfei Zhou, Fisher Yu, Jifeng Dai, Ender Konukoglu, and Luc Van Gool. 2021. Exploring cross-image pixel contrast for semantic segmentation. In *ICCV*.

[42] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. 2021. Dense contrastive learning for self-supervised visual pre-training. In *ICCV*.

[43] Yabin Wang, Zhiheng Ma, Xing Wei, Shuai Zheng, Yaowei Wang, and Xiaopeng Hong. 2022. Eccnas: Efficient crowd counting neural architecture search. *ACM TOMM* (2022).

[44] Fangyun Wei, Yue Gao, Zhirong Wu, Han Hu, and Stephen Lin. 2021. Aligning pretraining for detection via object-level contrastive learning. *NIPS* (2021).

[45] Xing Wei, Yuanrui Kang, Jihao Yang, Yunfeng Qiu, Dahu Shi, Wenming Tan, and Yihong Gong. 2021. Scene-Adaptive Attention Network for Crowd Counting. *arXiv preprint* (2021).

[46] Kilian Q Weinberger and Lawrence K Saul. 2009. Distance metric learning for large margin nearest neighbor classification. *Journal of machine learning research* (2009).

[47] Enze Xie, Jian Ding, Wenhai Wang, Xiaohang Zhan, Hang Xu, Peize Sun, Zhenguo Li, and Ping Luo. 2021. Detco: Unsupervised contrastive learning for object detection. In *ICCV*.

[48] Zhaoyi Yan, Yuchen Yuan, Wangmeng Zuo, Xiao Tan, Yezhen Wang, Shilei Wen, and Errui Ding. 2019. Perspective-guided convolution networks for crowd counting. In *ICCV*.

[49] Yifan Yang, Guorong Li, Zhe Wu, Li Su, Qingming Huang, and Nicu Sebe. 2020. Weakly-supervised crowd counting learns from sorting rather than locations. In *ECCV*.

[50] Lingke Zeng, Xiangmin Xu, Bolun Cai, Suo Qiu, and Tong Zhang. 2017. Multi-scale convolutional neural networks for crowd counting. In *ICIP*.

[51] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. 2016. Single-Image Crowd Counting via Multi-Column Convolutional Neural Network. In *CVPR*.

[52] Xiangyun Zhao, Raviteja Vemulapalli, Philip Andrew Mansfield, Boqing Gong, Bradley Green, Lior Shapira, and Ying Wu. 2021. Contrastive Learning for Label Efficient Semantic Segmentation. In *ICCV*.

| $\beta$ | 0 | 0.1 | 0.5 | 1 |
|---|---|---|---|---|
| MAE | 87.2 | 86.7 | 86.1 | **85.4** |
| MSE | 140.8 | 136.7 | 137.6 | **134.5** |
| $\beta$ | 2 | 5 | 10 | GP [32] |
| MAE | 87.1 | 87.2 | 87.0 | 102.0 |
| MSE | 134.8 | 136.2 | 138.9 | 172.0 |

Table 6: The influence of noise depression $\beta$. When $\beta$ is greater than $0$, the model accuracy can be improved. It indicates that reducing the weights of noisy supervisions is beneficial to the model.

| $\lambda_c$ | 0 | 0.001 | 0.01 | 0.05 |
|---|---|---|---|---|
| MAE | 94.2 | 85.5 | **85.4** | 85.7 |
| MSE | 148.9 | 135.1 | **134.5** | 134.6 |
| $\lambda_c$ | 0.1 | 0.5 | 1 | GP [32] |
| MAE | 88.9 | 89.9 | 92.6 | 102.0 |
| MSE | 136.4 | 135.9 | 139.9 | 172.0 |

Table 7: The influence of agency loss parameter $\lambda_c$. With the moderate help of the density agency, the accuracy of the counting model has been significantly improved.
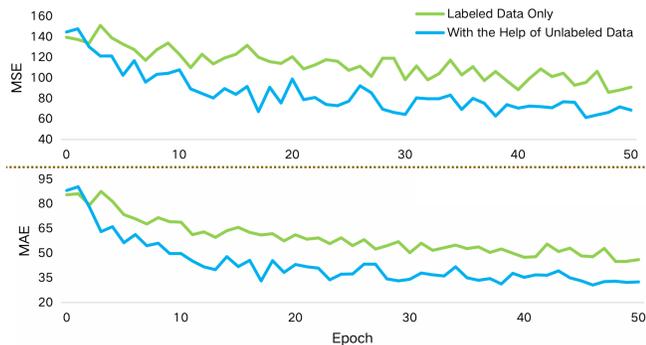


Figure 7: The curves of training MAE and MSE under two different settings: training with only labeled data and training with both labeled and unlabeled data. With the help of unlabeled data, the accuracy of the model improves faster and performs consistently better.

## A ABLATION STUDY

We hold more extensive experiments to study the influence of different parameters and settings. All experimental results are held in ShanghaiTech A with *a labeled ratio of* 5%. We choose previous state-of-the-art GP [32] for comparison. The final performances of our model under different settings are consistently better than the previous method.

### A.1 The convergence speed

To examine the impact of unlabeled data, we record the MAE and MSE of training epochs. The comparison result is shown in Figure 7, where the green curve denotes the model trained only by labeled data and the blue one denotes the model trained by both labeled and unlabeled data.

For the first 20 epochs, we find that with the help of unlabeled data, the error in training samples drops faster than that of training with only labeled data. And after that, this low training error is maintained and consistently better. This quantitative experiment proves that our model can get great help from unlabeled data to improve the performance.

### A.2 The influence of noise depression

We conduct experiments to study the influence of noise depression parameter $\beta$, which controls the strength of penalty. The comparison result is shown in Table 6.

When $\beta = 0$, the loss degenerates into normal Bayesian loss. And when we focus on $\beta > 0$, the results are constantly better than using all annotations indiscriminately. As $\beta$ is set larger, the depression of instances with large gaps will be more severe, but the final accuracy of the model becomes worse. It may suggest that when the depression strength is too strong, some appropriate supervision signals are also ignored, which leads to a deterioration of performance.

### A.3 The influence of density agency

The parameter $\lambda_c$ determines how much that the density agency will supervise the model. Specifically, when $\lambda_c = 0$, the density agency will not work and the model is only supervised by ground-truth of labeled data. The experimental result is shown in Table 7.

As it shown, without the help of density agency, the performance of the counting model drops a lot. And when $\lambda_c = 0.01$, the supervisions between density agency and the ground-truth achieves a good balance, where the model performs the highest accuracy. However, as $\lambda_c$ continues to increase, the proportion of this supervision becomes larger, gradually replacing the supervisions from ground-truth. The accuracy of the model gradually decreases, even approaching the performance of not using the agency.

### A.4 The influence of temperature parameter

We conduct experiments of the temperature parameter $\tau$ in density agent guided contrastive learning. When $\tau$ is larger, the difference between different density features is smaller. The supervision of contrastive learning will be weakened. The result is shown in Table 8.

In general, when tuning $\tau$, the final performance does not fluctuate much and is consistently quite better than the previous method. Particularly, when $\tau$ is too large or too small, the accuracy will still be affected to a certain extent. During other experiments, we keep $\tau = 0.1$.

### A.5 The influence of mask loss

The mask loss parameter $\lambda_m$ controls the weight of the supervision targeted on foreground predictor. We conduct experiments to study the influence of mask loss, which result is shown in Table 9.

Since dividing foreground and background is important for subsequent transformer and other structures, setting a $\lambda_m$ that is too small is inappropriate. As the experiment shows, when $\lambda_m = 0.01$,

| $\tau$ | 0.01 | 0.05 | 0.07 | 0.1 |
|---|---|---|---|---|
| MAE | 89.6 | 87.8 | 86.0 | **85.4** |
| MSE | 138.4 | 136.3 | 135.8 | **134.5** |

| $\tau$ | 0.2 | 0.5 | 1 | GP [32] |
|---|---|---|---|---|
| MAE | 85.4 | 86.7 | 88.9 | 102.0 |
| MSE | 134.7 | 135.6 | 139.2 | 172.0 |

**Table 8: The influence of temperature parameter. In general, the adjustment of contrastive temperature has little fluctuation on the final result.**

| $\lambda_m$ | 0.01 | 0.05 | 0.01 | 0.5 | 1 |
|---|---|---|---|---|---|
| MAE | 90.8 | 87.9 | **85.4** | 85.8 | 90.5 |
| MSE | 142.0 | 137.7 | 134.5 | **134.4** | 146.7 |

**Table 9: The influence of mask loss. When the loss weight $\lambda_m$ is too small or too large, the performance of the counting model will deteriorate.**

the performance of the counting model deteriorates. However, when the weight is too large, that is $\lambda_m = 1$, the accuracy also drops. It can be explained as the model pays too much attention to the division of the foreground and background, while ignoring the difference between different foreground densities.