# Coarse- and Fine-grained Attention Network with Background-aware Loss for Crowd Density Map Estimation

Liangzi Rong      Chunping Li
School of Software, Tsinghua University
lzrong13@gmail.com      cli@tsinghua.edu.cn

## Abstract

*In this paper, we present a novel method Coarse- and Fine-grained Attention Network (CFANet) for generating high-quality crowd density maps and people count estimation by incorporating attention maps to better focus on the crowd area. We devise a from-coarse-to-fine progressive attention mechanism by integrating Crowd Region Recognizer (CRR) and Density Level Estimator (DLE) branch, which can suppress the influence of irrelevant background and assign attention weights according to the crowd density levels, because generating accurate fine-grained attention maps directly is normally difficult. We also employ a multi-level supervision mechanism to assist the backpropagation of gradient and reduce overfitting. Besides, we propose a Background-aware Structural Loss (BSL) to reduce the false recognition ratio while improving the structural similarity to groundtruth. Extensive experiments on commonly used datasets show that our method can not only outperform previous state-of-the-art methods in terms of count accuracy but also improve the image quality of density maps as well as reduce the false recognition ratio.*

## 1. Introduction

Recently, crowd density map estimation has received continuous attention as a challenging computer vision task. Given a crowd image, its purpose is to estimate the density map and the total number of people. The value of each pixel in the density map reflects the density of the corresponding area in the image, and estimated people count can be obtained by accumulating the values of all pixels. Generally speaking, there are three major difficulties for accurate crowd counting: (1) Different distances from the shooting device lead to scale variation within one image and between different images. (2) Severe occlusion in high-density crowd scenes. (3) The influence of complex and irrelevant background is not conducive to recognizing crowd areas. Existing methods[1, 28, 11, 13] mainly adopt multi-
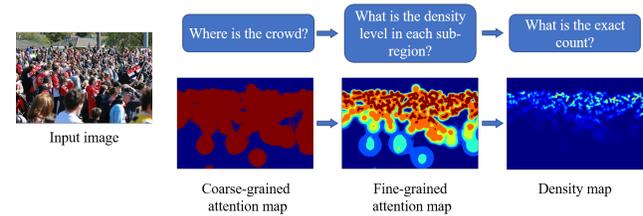


Figure 1. An intuitive explanation of our method.

scale or multi-column architecture-based CNN models to capture richer features. But they generally have two disadvantages: (1) The influence of the background is not fully considered. Existing methods treat all the regions in one image as potential crowd areas, so the convolution kernels may extract features in the background regions that are not related to the crowd, resulting in misrecognition. In this case, even if the estimated number of people across the whole image is close to the ground truth, it may be caused by both the underestimation of the crowd region and the misrecognition of the background region. (2) All crowd regions are treated equally. In fact, when a person observes a crowd image, it is normal to pay different attention to different areas. This is because for areas with a high degree of occlusion, it is more difficult to distinguish the features of each person, thus worth more attention. In low-density areas, it is easier to distinguish each person, thus worth less attention. However, most existing methods pay the same attention to all sub-regions with in one image.

In order to solve the above two problems, we hope to design an attention mechanism that can suppress the influence of the background and reduce misrecognition, but also adaptively assign attention weights to regions with different density. Therefore, we introduce the attention maps. We expect that in one attention map, the background area has a weight close to 0, and the low-density area has relatively low weights, and the high-density area has relatively high weights. As depicted in Fig. 1, when the human eyes observe a picture, they will first recognize in which areas people exist, and then identify the dense degree and count.

Intuitively, it is natural to follow this paradigm when designing a CNN. It is easier to judge whether there are people in an area than to judge the density level of the area, so we can get more reliable results if we start from simple tasks. Therefore, we employ two modules, Crowd Region Recognizer (CRR) and Density Level Estimator (DLE), to judge whether people exist and the density level for an area respectively.

To be more specific, CRR produces a coarse-grained attention map (CAM), whose value indicating the possibility of if people exist. DLE produces a fine-grained attention map (FAM), whose value indicating the density level of each area. In order to leverage the relatively reliable CAM, we combine the FAM with it. Then feature maps for regressing density maps are combined with FAM on multiple scale to pay adaptive attention to different areas.

The images of crowd datasets is relatively few compared to other datasets, e.g. Imagenet. To reduce overfitting and facilitate the backpropagation of gradients, we introduce multi-level supervision by adding several more output layers in internal layers and summing all the loss functions and backpropagate.

In addition, this paper also explores the impact of different loss functions. We propose a novel and effective loss function named Background-aware Structural Loss (BSL) that can take account of structural similarity, counting accuracy and false recognition ratio. We evaluate the performance on multiple datasets and models, and proposed BSL delivers superior performance than other loss functions.

To sum up, our contribution are three-fold: (1) We present Coarse- and Fine-grained Attention Network (CFANet) that can produce high-quality density maps and accurate count estimation. Crowd Region Recognizer (CRR) and Density Level Estimator (DLE) are employed to estimated coarse- and fine-grained attention maps respectively to help pay more attention to areas with high density and less attention to areas without people or with a low density. (2) We introduce the multi-level supervision mechanism to facilitate the backpropagation of gradients and reduce overfitting. (3) We propose a loss function named Background-aware Structural Loss (BSL) that can improve structural similarity and reduce false recognition. By combining CFANet with BSL, we can achieve the best performance on most mainly used datasets.

## 2. Related Work

We classify existing methods into two categories, one is only using the density map as the learning objective, and the other is combining classification, segmentation and other tasks as the learning objectives.

**Density Map as the Learning Objective.** To solve the problem that the size of the heads in the crowd scene changes greatly, MCNN[29] designed three paral-

lel subnetworks with convolution kernels of different sizes. CSRNet[8] found that the subnetworks introduced redundant parameters, so a single-column model was proposed. SANet[1] combined convolution kernels of different sizes as the feature extraction part, which enhances the non-linear expression ability of the network. CAN[14] used pooling pyramid to extract features of different scales, and adaptively assigns different weights to different scales and regions. TEDNet[7] proposed a lightweight hierarchical network structure by combining high-level and low-level features effectively. RR[21] used representative images in the dataset obtained by clustering to assist regress the residual and final density map. Bayesian loss[15] was proposed based on the Bayesian probability model. The value of each pixel in the density map contributes to each person's value or the background by probability. L2SM[26] further zooms in and re-estimates the dense area to improve the performance in the high-density area after obtaining the initial density map. DSSINet[11] uses conditional random fields to enable the features of each scale to obtain information from other scales for improvement.

**Combinal Learning Objectives.** CP-CNN[19] uses extra classification networks to combine global and local density level to improve model accuracy. Switch-CNN[17] classifies the input images into three categories according to density level, and then processes the images with corresponding sub-network. DecideNet[10] constructed a model using two sub-networks for detection and regression respectively, and the outputs of the two sub-networks are fused using attention mechanism. RAZNet[9] uses an additional localization branch to detect the heads' position, and adaptively zooms in the regions that are difficult to recognize. SGANet[22] employs the Inception-V3 structure as the feature extraction part, and uses the penultimate layer's feature maps to regress an attention map.

The most similar work to ours is ADCrowdNet[13] which uses a separate classification network to obtain the attention map, and multiplies it with the image or feature maps. We differ from it from three aspects: (1) We find that single scale combination has limited influence so we combine the attention maps on multiple scales to enhance the attention effect. (2) Attention maps in ADCrowdNet is similar to our coarse attention maps, but we further incorporate the density level information, which is not included in ADCrowdNet. (3) We do not need a separate classification network to produce attention maps, so our method is much easier for end-to-end training.

## 3. Our Approach

The architecture of the proposed method is illustrated in Figure 2. It consists of four modules: Feature map extractor is used to extract general feature maps from crowd images and feed them into the following modules for fur-
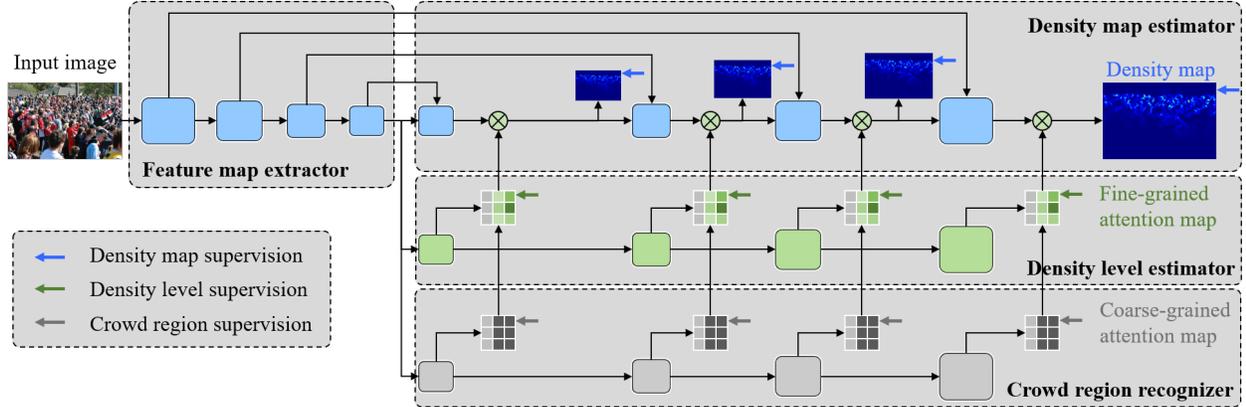
Figure 2. The structure of proposed CFANet.

ther processing. Crowd region recognizer (CRR) is used to judge whether an area contains crowd, and outputs a coarse-grained attention map (CAM). Density level estimator (DLE) is used to further classify the density level of each area, and can output a fine-grained attention map (FAM). With the aid of FAM, density map estimator can better focus on the crowd region and produce high-quality density maps. In addition, to assist the backpropagation and reduce overfitting, we design a multi-level supervision mechanism. In each stage of the density map estimator, feature maps are upsampled and fed into $Conv$ layers to obtain a density map, and the loss functions of multiple stages are summed and backpropagated.

## 3.1. Feature map extractor

Coarse- and fine-grained attention map and density map, which will be detailed in the following subsections, share a feature: having high values in the high-density area, and low values in the low-density area and near zero in the background area. Therefore, we argue that general features can be extracted through the same base network.

In most existing studies, the resolution of estimated density maps are lower than input, leading to the loss of spatial details. Inspired by the success of UNet[16], we design our model in an encoder-decoder paradigm. Following the practice of most previous work, we adopt the VGG-16's feature extraction part in this part. As Fig. 2 shows, we retain the first 10 convolutional layers and 3 pooling layers. Therefore, we can get feature maps with sizes of 1, 1/2, 1/4, and 1/8 from each stage.

## 3.2. Crowd region recognizer (CRR)

Because the crowd images contain different scenes, accurate crowd counting may be hindered by complex backgrounds. Even if the overall estimated number of people is close to the groundtruth, it may be caused by the underestimation of the crowd area and the false recognition of the background area. To address this issue, estimated coarse-grained attention maps (CAM) from CRR try to classify each pixel in feature maps into two categories: crowd and background region. Detailed configuration of the CRR module is: C(256,3)-U-C(128, 3)-U-C(128, 3)-U-C(64, 3)-C(1, 3), where C means convolution layer and U means bilinear upsample layer with rate=2. At each stage in CRR, the feature maps are fed into a $3 \times 3$ $Conv$ layer to regress a coarse-grained attention map which is then fed into corresponding stage in DLE. Losses calculated on multi-stage are summed and backpropagated as the gray arrow in Fig. 2. We will describe how to generate the groundtruth CAM, i.e., learning objective of this module, in Section 4.1.

## 3.3. Density level estimator (DLE)

The CRR module only implements a coarse-grained attention mechanism, that is, it only distinguishes between people and background. The goal of the DLE module is to further classify the crowd area into different density levels, because for high-density areas, we should give more attention, and for low-density areas, we should give less attention. Similar to CRR, this module should also try to suppress the influence of unmanned background. To this end, we divide all pixels into $k$ categories according to the threshold obtained from statistics and consider it as a $k$-class classification problem. This process of the generating groundtruth FAM, i.e., learning objective of DLE module, is detailed in Section 4.1. After the $k$-class classification, different attention weights are assigned accordingly: pixels classified to class 0 are regarded as the background, and the attention weight are set to 0. For other density levels, we divide the range of (0,1] into k-1 categories to correspond. For example, when k is 6, classes 1-5 correspond to 0.2, 0.4, 0.6, 0.8 and 1 respectively. Regarding the number of categories, we have conducted an experimental comparison, and found that when the number of categories is set to 6, the performance is best. Detailed experimental results can be
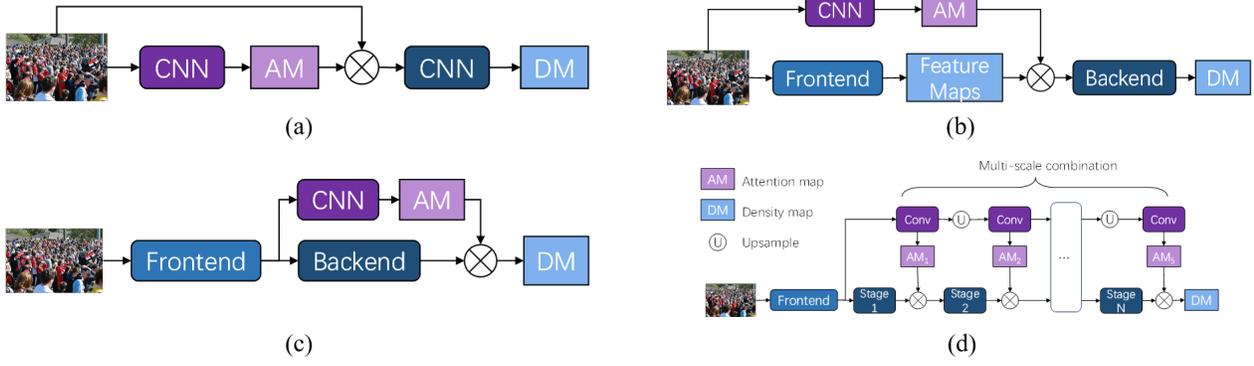
Figure 3. Comparison of attention-based methods. AME means attention map estimator. DME means density maps estimator. Proposed method is (d).

seen in section 4. Detailed configuration of the DLE module is: C(256,3)-U-C(256, 3)-U-C(128, 3)-U-C(64, 3)-C($k$, 3), where C means convolution layer and U means bilinear upsample layer. At each stage in DLE, feature maps are fed into a $3\times3$ $Conv$ layer and regress a fine-grained attention map. Similar to CRR, losses calculated on multi-stage are summed and backpropagated as the green arrow in Fig. 2. Since we already have the CAM from CRR, we can make full use of it to refine the FAM:

$$FAM = FAM + CAM \qquad (1)$$

Then the FAM is fed into the corresponding stage in the next part.

### 3.4. Density map estimator

This module's configuration is almost symmetrical with the feature map extractor. Since CFANet aims to estimate high-resolution and high-quality density maps, DME is constructed with a set of convolutional and upsample layers: C(512, 3)-U-C(256, 3)-U-C(256, 3)-U-C(64, 3)-C(1, 1), where C is convolution and U bilinear upsample layer. We set dilation rate = 2 in all convolutional layers to enlarge the receptive field.

There have been studies using the attention map to better focus on the crowd area. As Fig. 3 shows, (a) and (b) train a separate network to predict the attention map. (a) multiplies the input image with predicted attention map and then use it for regression. (b) multiplies feature maps extracted by the frontend with the predicted attention map. (c) inserts an attention branch to predict attention map before the last convolutional layer. The density map is multiplied with the predicted attention map as the final output. It is worth noting, however, that the weights assigned to the background area by estimated attention map is not equal to 0, so the influence cannot be completely suppressed with one-time combination. Therefore, we adopt a multi-scale attention combination method, as shown in Figure 3 (d). At each

stage in this module, the feature maps (FM) are combined with the fine-grained attention maps (FAM) from the DLE in a residual way:

$$FM = FM + FAM * FM \qquad (2)$$

In this way, the attention effects are enhanced and different areas are paid adaptive attention with the aid of fine-grained attention maps.

We also employ a multi-level supervision mechanism. Feature maps in each stage will be upsampled to the size of the input image, and fed into a $3\times3$ convolutional layer to regress a density map as the blue arrow in Fig. 2 indicates. Therefore, for each input image, CFANet will produce 4 density maps and the 4 losses of are summed and backpropagated. Generally, a deeper network has stronger expression ability, so we regard the output of the last layer as the final output.

### 3.5. Loss function

We first introduce a structural loss function (SL) that considers both structural similarity and counting accuracy. The definition is as follows:

$$SL = \frac{1}{K}\sum_{j=1}^{K}(1 - SSIM(Pool_j(DM), Pool_j(\hat{DM})))$$

$$(3)$$

$$SSIM(X,Y) = 1 - (\frac{(2\mu_X\mu_Y + C_1)(2\sigma_{XY} + C_2)}{(\mu_X^2 + \mu_Y^2 + C_1)(\sigma_X^2 + \sigma_Y^2 + C_2)})$$

$$(4)$$

where $DM$ and $\hat{DM}$ mean groundtruth and estimated density map and $Pool_j$ means downsampling to $\frac{1}{2^{j-1}}$ size with average pooling. $\mu$ denotes the local mean and $\sigma$ is the local variance, $\sigma_{XY}$ is the local covariance. $C_1$ and $C_2$ are set to 0.01 and 0.03. $K$ is set to 3. SSIM of high resolution density maps can focus on spatial details and improve structural similarity. SSIM of downsampled density maps can improve global counting accuracy.

To reduce the false recognition ratio, we add a background-aware loss item (BL):

$$BL = \frac{C_{bg}}{C_{total}} \qquad (5)$$

where $C_{total}$ is the estimated total people count and $C_{bg}$ is the estimated people count in the background area. The background area is divided in the same way as the groundtruth CAM which is detailed in section 4.1.

The loss function (BSL) for density map optimization is the sum of SL and BL.

To optimize the coarse- and fine-grained attention maps, we use cross-entropy as the loss function.

The final loss function is the weighted sum of density map and attention maps loss functions at multiple stages.

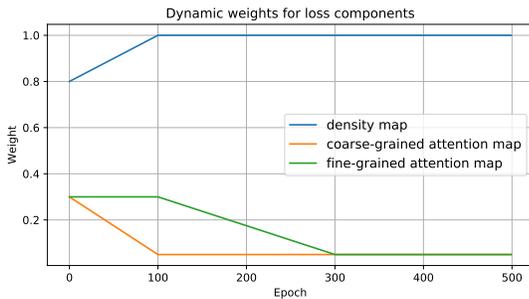$$L_{total} = SL + BL + \lambda L_{CAM} + \mu L_{FAM} \qquad (6)$$



Figure 4. Dynamic weight adjustment strategy.

Since it is the easiest task to divide the image into background and crowd, and the density multi-classification has medium difficulty, and it is the most difficult to regress the accurate density value of each pixel, we adopt a dynamic weight adjusting strategy in the loss function. In the initial training stage, the network cannot accurately capture the distribution features. Therefore, we focus more on the easy tasks, and set the weights of binary classification and multi-class classification large. As the training process progresses, the representative ability of the network continues to improve, and the features can be better captured. At this time, the weight of density map loss is adjusted to be larger. This process is shown in the Fig. 4.

# 4. Experiments and Analysis

## 4.1. Datasets and Evaluation Criteria

We evaluate the performance of proposed CFANet on four major crowd counting datasets: ShanghaiTech [29], UCF_CC_50 [5], UCF-QNRF [6] and Mall [2]. The detailed comparisons are described in the following subsections.

**Groundtruth Generation.** We first generate groundtruth density maps and then generate groundtruth

coarse- and fine-grained attention maps. For density maps, following [29], each labeled head $p_j$ is substituted with a Gaussian kernel $\mathcal{N}(p_i, \sigma^2)$, where $\sigma$ is the average distance of $p_j$ and its 3 nearest neighbors. The kernel is normalized to 1, so the integral of the density map is equal to the labeled people count. For CAM, we set the value to 1 for one pixel if the value $\geq 1e - 5$ for corresponding position in density map and 0 otherwise to obtain the groundtruth. For FAM, we categorize one pixel to class 0 if its value $< 1e - 5$ in density map, and then count all the other values and divide them into $k - 1$ categories from large to small, and each pixel's class is assigned accordingly.

We adopt Mean Absolute Error(MAE) and Root Mean Square Error(RMSE) for evaluating the counting accuracy:

$$MAE = \frac{1}{N} \sum_{i=1}^{N} \left| Count_{est}^i - Count_{gt}^i \right| \qquad (7)$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (Count_{est}^i - Count_{gt}^i)^2} \qquad (8)$$

SSIM and PSNR are used for evaluating the quality of density maps.

## 4.2. Implementation Details

During training, each image is random cropped to $\frac{1}{2}$ size, then horizontal flipped at possibility 0.5 to enlarge the dataset. Pre-trained VGG-16 on ImageNet is used to initialize the parameters of the feature map extractor, and other parameters are randomly initialized by Gaussian distribution with $\delta = 0.01$. We train our network with Adam optimizer for 500 epochs and learning rate is initially set to 2e-5 and reduced by half every 100 epochs. We also multiply the density maps by an expansion factor=50 as the directly generated pixel values are quite small.

## 4.3. Comparison with State-of-the-art Methods

We compare our model with state-of-the-art crowd density maps estimation methods since 2016.

**ShanghaiTech** dataset consists of PartA and PartB. PartA has 300 training images and 182 testing images with relatively high density. PartB has 400 training images and 316 testing images with relatively low density. As Table 1 shows, CFANet outperforms the state-of-the-art methods. Specifically, we improve the MAE by 1.6% and achieve the second best RMSE on PartA sub-dataset, and improve the MAE by 3.0% and the RMSE by 1.0% on PartB sub-dataset. These results suggest that our method can be well applied in both crowded and sparse scenes.

**UCF_CC_50** dataset includes 50 crowd images with extremely high density. This dataset is quite challenging not

| Method | PartA | | PartB | |
|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE |
| CSRNet[8] | 68.2 | 115.0 | 10.6 | 16.0 |
| SANet[1] | 67.0 | 104.5 | 8.4 | 13.6 |
| TEDNet[7] | 64.2 | 109.1 | 8.2 | 12.8 |
| ADCrowdNet[13] | 63.2 | 98.9 | 7.7 | 12.9 |
| AT-CSRNet[30] | - | - | 8.1 | 13.5 |
| CAN[14] | 62.3 | 100.0 | 7.8 | 12.2 |
| Bayesian loss[15] | 62.8 | 101.8 | 7.7 | 12.7 |
| DSSINet[11] | 60.6 | 96.0 | 6.9 | <u>10.3</u> |
| RANet[28] | 59.4 | 102.0 | 7.9 | 12.9 |
| S-DCNet[25] | 58.3 | 95.0 | <u>6.7</u> | 10.7 |
| PGCNet[27] | <u>57.0</u> | **86.0** | 8.8 | 13.7 |
| **Ours** | **56.1** | <u>89.6</u> | **6.5** | **10.2** |

Table 1. Comparison of performance on ShanghaiTech dataset. Top two performance are highlighted in **bold** and <u>underline</u>.

only because of the quite limited labeled images, but also because the people count in each image varies from 94 to 4543 with an average of 1279.5. We follow the standard 5-fold cross-validation in [5]. As shown in Table 2, we achieve the best MAE and the second best RMSE performance. Although this dataset has very limited labeled images, proposed CFANet can still get a promising result, which shows its capacity when there is a training sample shortage.

| Method | MAE | RMSE |
|---|---|---|
| CSRNet[8] | 266.1 | 397.5 |
| ADCrowdNet[13] | 257.1 | 363.5 |
| CAN[14] | 212.2 | **243.7** |
| PGCNet[27] | 244.6 | 361.2 |
| Bayesian loss[15] | 229.3 | 308.2 |
| SANet[1] | 258.4 | 334.9 |
| TEDNet[7] | 249.4 | 354.5 |
| DSSINet[11] | 216.9 | 302.4 |
| RANet[28] | 239.8 | 319.4 |
| S-DCNet[25] | <u>204.2</u> | 301.3 |
| **Ours** | **203.6** | <u>287.3</u> |

Table 2. Comparison of performance on UCF_CC_50 dataset. Top two performance are highlighted in **bold** and <u>underline</u>.

**UCF-QNRF** is the most up-to-date large scale crowd image dataset. It has 1535 labeled images in total, with 1201 images for training and 334 images for testing. The people count in each image ranges from 49 to 12865, which makes it quite challenging and diverse. Because the resolutions vary in a large range, we resize the images to ensure that the longer side falls in [1024, 2048] in the training process. In the testing process, we use the original images.

The comparison is summarized in Table 3, our method delivers the best RMSE, surpassing the second best approach by 1.6%, and the second best MAE, which is quite close to the best one.

| Method | MAE | RMSE |
|---|---|---|
| CAN[14] | 107 | 183 |
| TEDNet[7] | 113 | 188 |
| RANet[28] | 111 | 190 |
| S-DCNet[25] | 104.4 | 176.1 |
| SFCN[23] | 102.0 | 171.4 |
| DSSINet[11] | 99.1 | 159.2 |
| MBTTBF-SFCB[20] | 97.5 | 165.2 |
| Bayesian loss[15] | **88.7** | <u>154.8</u> |
| **Ours** | <u>89.0</u> | **152.3** |

Table 3. Comparison of performance on UCF-QNRF dataset. Top two performance are highlighted in **bold** and <u>underline</u>.

**Mall** dataset [2] consists of 2000 frames obtained from a stationary surveillance camera in a mall. The images have a resolution of 320×480. Following the same setting as [2], the first 800 frames are used as training frames and the remaining 1200 frames are used for testing. It can be viewed from Table 4 that our method get the best performance in both MAE and RMSE, improving the performance by 6.3% and 7.1%, which demonstrates that our method also has the superior capacity in low-density scenes.

| Method | MAE | RMSE |
|---|---|---|
| ConvLSTM[24] | 2.10 | 7.6 |
| DRSAN[12] | 1.73 | 2.1 |
| DecideNet[10] | 1.52 | 1.90 |
| AT-CFCN[30] | 2.28 | 2.90 |
| E3D[31] | 1.64 | 2.13 |
| SAAN[4] | <u>1.28</u> | <u>1.68</u> |
| **Ours** | **1.20** | **1.56** |

Table 4. Comparison of performance on Mall dataset. Top two performance are highlighted in **bold** and <u>underline</u>.

We select representative images from each dataset and compare the predicted density map with groundtruth in Fig. 5. It can be viewed that the distribution of crowd is very close to the groundtruth.



Figure 5. Visualization of density maps. Three rows from top to bottom are input image, groundtruth and estimated density maps from CFANet respectively.

Some feature maps from internal layers are visualized in

6. Refined feature maps in the last column retain most feature in crowd areas and discard unrelated features in background areas and help reduce false recognition ratio.
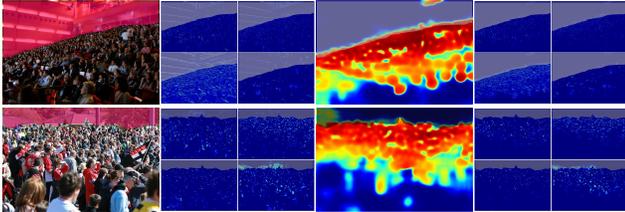


Figure 6. Feature maps before and after refined by fine-grained attention maps. Four columns from left to right are input image, original feature maps, fine-grained attention map (FAM) and feature maps refined by FAM. The background area is highlighted by the mask.

We also compare the density map quality under SSIM and PSNR criterion. The performances are summarized in Table 5. On the ShanghaiTech PartA dataset, density maps produced by our method have the best quality. Note that our estimated density maps have the same resolution as the input image. This indicates that proposed CFANet has the capacity of improving the quality of high-resolution density maps as well as predicting accurate people count.

| Method | SSIM | PSNR |
|---|---|---|
| MCNN[29] | 0.52 | 21.40 |
| CSRNet[8] | 0.76 | 23.79 |
| ADCrowdNet[13] | **0.88** | 24.48 |
| TEDNet[7] | 0.83 | 25.88 |
| **Ours** | **0.88** | **30.11** |

Table 5. SSIM and PSNR comparisons on ShanghaiTech PartA dataset.

As described previously, we propose this method to suppress the influence of irrelevant background and reduce false recognition. To find out if CFANet can function well in this aspect, we set a criterion for judging the ratio of false recognition in the background. In the testing process, we separate the estimated people count in the background and crowd area, and calculate the proportion of these two parts to the total number of people. To be more specific, for each image in the testing set, the pixels with value$< 1e - 5$ are considered as in the background area, and others are considered as in the crowd area. By accumulating the pixel values of the background area, we can obtain the corresponding estimated number of people: $Count_{est}^{bg}$. The ratio $r_{bg} = Count_{est}^{bg}/Count_{est}$ represents the false recognition ratio, and the lower it is, the less false recognition there is. We compare the $r_{bg}$ of proposed method with two open-source methods CSRNet[1] and CAN[2] and removing the BL

---

[1] https://github.com/leeyeehoo/CSRNet-pytorch
[2] https://github.com/weizheliu/Context-Aware-Crowd-Counting

---

item (Equation 5) in loss function. Results in Table 6 validate the effectiveness of CFANet and background-aware loss item in suppressing the false recognition in the background area.

| Method | PartA | PartB |
|---|---|---|
| CSRNet[8] | 0.0381 | 0.0857 |
| CAN[14] | 0.0283 | 0.0359 |
| Ours w/o bg-aware loss | 0.0205 | 0.0245 |
| Ours w. bg-aware loss | **0.0184** | **0.0219** |

Table 6. Comparison of false recognition ratio in estimated density maps on ShanghaiTech dataset.

## 4.4. Ablation Study

We conduct extensive ablation studies on ShanghaiTech PartA dataset to analyze the impact of different settings of the network and training process.

**The effectiveness of CRR and DLE**. Since we have introduced two new branches: CRR and DLE, we first try to analyze the impact of branch setting and identify if they can bring improvement to the performance. We construct the 'baseline' by removing the CRR and DLE. We then add the CRR and DLE respectively to build 'w. CRR' and 'w. DLE'. Finally, both branches are integrated as 'w. CRR+DLE', which is the default setting in experiments. The results are reported in Table 7. We can see that the integration of CRR and DLE branch can both bring improvement, demonstrating the effectiveness of attention mechanism, and by combining them the best performance can be achieved.

| Network | MAE | RMSE |
|---|---|---|
| baseline | 63.7 | 102.9 |
| w. CRR | 56.9 | 91.8 |
| w. DLE | 60.2 | 99.4 |
| w. CRR+DLE | **56.1** | **89.6** |

Table 7. Ablation study on CRR and DLE branch configuration.

**Multi-level supervision**. Table 8 reports the performance of using different supervision's setting. $Supervision\ i$ corresponds using the $i$-th blue arrow in Fig. 2. We can see that the commonly used setting, single $Supervision\ 4$, can only produce a relatively low accuracy. The more supervision is integrated, the more improvement it can bring to the model's accuracy and using all 4 supervision achieves the best performance. This undoubtedly demonstrates the superiority of the multi-level supervision mechanism.

**Class number in DLE branch**. As described previously, we need to specify a class number for the multi-class classification problem in DLE branch. If the class number is too small, the features of each category are not representative enough. If the number is too big, the differences

| Supervision | MAE | RMSE |
|---|---|---|
| $Supervision\ 4$ | 69.5 | 109.1 |
| $Supervision\ 3 \sim 4$ | 64.6 | 106.3 |
| $Supervision\ 2 \sim 4$ | 62.0 | 97.1 |
| $Supervision\ 1 \sim 4$ | **56.1** | **89.6** |

Table 8. Ablation study on supervision setting.

between classes can be too little to capture. To find out an appropriate setting, we have tried setting it to 4, 6, 8, 10, and experimental results in Table 9 show that setting it to 6 can achieve the best performance.

| Class number | MAE | RMSE |
|---|---|---|
| 4 | 56.9 | 91.3 |
| 6 | **56.1** | **89.6** |
| 8 | 57.4 | 91.7 |
| 10 | 57.9 | 92.9 |

Table 9. Ablation study on density level class number setting.

# 5. The Impact of Loss Functions

Many kinds of loss functions have been proposed to improve the performance of density map estimation models. Here we firstly describe existing loss functions briefly, and then compare proposed BSL with them.

Mean Square Error ($MSE$) is the most widely used loss function which computes the average of the square error of each pixel in density maps.

Spatial Abstraction Loss ($L_{SA}$) and Spatial Correlation Loss ($L_{SC}$) are proposed by [7]. $L_{SA}$ computes $MSE$ losses on multiple abstraction levels and $L_{SC}$ further computes global consistency.

Mean Absolute Error ($MAE$) is introduced by [18] to add robustness to outliers. $MAE$ and $MSE$ are jointly optimized in training.

Multi-scale density level consistency loss ($L_c$) is proposed by [3] which combines pixel-wise $MSE$ with $MAE$ of globally averaged density maps.

Bayesian loss ($L^{Bayes+}$) is introduced by [15], which constructs a density contribution probability model from the point annotations, instead of constraining the value at every pixel in the density map.

SSIM is firstly introduced by [1] as loss function ($L_{SSIM}$) to improve the quality of results by incorporating the local correlation in density maps.

Dilated Multiscale Structural Similarity (DMS-SSIM) loss ($L_{DS}$) is introduced by [11]. The DMS-SSIM network for computing the $L_{DS}$ consists of $m = 5$ dilated convolutional layers with dilation rates of 1, 2, 3, 6 and 9.

## 5.1. Impact on Counting Accuracy

We evaluate the impact of loss functions on counting accuracy by training CFANet with the above existing loss functions and proposed BSL on ShanghaiTech PartA and UCF-QNRF datasets. As Table 10 shows, BSL achieves the best performance on both datasets, which validate its effectiveness.

| Loss function | ST PartA | | UCF-QNRF | |
|---|---|---|---|---|
| | MAE | RMSE | MAE | RMSE |
| $MSE$ | 63.6 | 104.5 | 101.9 | 170.4 |
| $L_{SA} + L_{SC}$[7] | 64.2 | 109.1 | 113.0 | 188.0 |
| $L_C$[3] | 65.2 | 107.7 | 101.3 | 174.2 |
| $MSE + MAE$[18] | 59.7 | 98.2 | 99.4 | 164.8 |
| $L_{Bayes+}$[15] | 66.4 | 113.9 | 97.2 | 160.6 |
| $L_{SSIM}$[1] | 64.9 | 107.4 | 108.8 | 179.5 |
| $L_{DS}$[11] | 60.6 | 96.0 | 99.1 | 159.2 |
| $BSL$ | **56.1** | **89.6** | **89.0** | **152.3** |

Table 10. Comparison of loss functions on ShanghaiTech PartA and UCF-QNRF datasets.

In order to further validate BSL's robustness on different models, we replace original $MSE$ loss function with proposed BSL on MCNN[29], CSRNet[8] and CAN[14] on ShanghaiTech PartA dataset. Results in Table 11 show that, without changing the network, using BSL can improve the performance on multiple models significantly.

| Model | Using MSE/BSL | |
|---|---|---|
| | MAE | RMSE |
| MCNN[29] | 110.2/**88.2**↑$^{20.0\%}$ | 173.2/**140.5**↑$^{18.9\%}$ |
| CSRNet[8] | 68.2/**63.1**↑$^{7.5\%}$ | 115.0/**102.5**↑$^{10.9\%}$ |
| CAN[14] | 62.3/**59.0**↑$^{5.3\%}$ | 100.0/**90.1**↑$^{9.9\%}$ |

Table 11. Comparison of BSL and MSE on multiple models.

# 6. Conclusion

In this paper, we present a new model named Coarse- and Fine-grained Attention Network (CFANet) for high-quality crowd density map generation and people counting. By incorporating Crowd Region Recognizer (CRR) and Density Level Estimator (DLE) to estimate coarse- and fine-grained attention maps, feature maps for regressing the density maps are refined on multi-scale and the network can better focus on the crowd region. We also adopt multi-level supervision to help facilitate the backpropagation of gradient and reduce overfitting. In addition, we propose a novel and effective loss function named Background-aware Structural Loss (BSL), which can achieve better counting accuracy and enhance structural similarity as well as reduce false recognition ratio. Combining proposed CFANet with BSL can outperform current state-of-the-art methods on most mainly used datasets.

# References

[1] 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 (ECCV)*, pages 734–750, 2018.

[2] Ke Chen, Chen Change Loy, Shaogang Gong, and Tony Xiang. Feature mining for localised crowd counting. In *BMVC*, volume 1, page 3, 2012.

[3] Feng Dai, Hao Liu, Yike Ma, Juan Cao, Qiang Zhao, and Yongdong Zhang. Dense scale network for crowd counting. *arXiv preprint arXiv:1906.09707*, 2019.

[4] Mohammad Hossain, Mehrdad Hosseinzadeh, Omit Chanda, and Yang Wang. Crowd counting using scale-aware attention networks. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1280–1288. IEEE, 2019.

[5] 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*, pages 2547–2554, 2013.

[6] Haroon Idrees, Muhmmad 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 (ECCV)*, pages 532–546, 2018.

[7] 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*, pages 6133–6142, 2019.

[8] 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*, pages 1091–1100, 2018.

[9] Chenchen Liu, Xinyu Weng, and Yadong Mu. Recurrent attentive zooming for joint crowd counting and precise localization. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1217–1226. IEEE, 2019.

[10] 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*, pages 5197–5206, 2018.

[11] 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 IEEE International Conference on Computer Vision*, pages 1774–1783, 2019.

[12] Lingbo Liu, Hongjun Wang, Guanbin Li, Wanli Ouyang, and Liang Lin. Crowd counting using deep recurrent spatial-aware network. *arXiv preprint arXiv:1807.00601*, 2018.

[13] 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*, pages 3225–3234, 2019.

[14] Weizhe Liu, Mathieu Salzmann, and Pascal Fua. Context-aware crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5099–5108, 2019.

[15] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. Bayesian loss for crowd count estimation with point supervision. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6142–6151, 2019.

[16] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[17] Deepak Babu Sam, Shiv Surya, and R Venkatesh Babu. Switching convolutional neural network for crowd counting. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4031–4039. IEEE, 2017.

[18] Zenglin Shi, Pascal Mettes, and Cees GM Snoek. Counting with focus for free. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4200–4209, 2019.

[19] Vishwanath A Sindagi and Vishal M Patel. Generating high-quality crowd density maps using contextual pyramid cnns. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1861–1870, 2017.

[20] Vishwanath A Sindagi and Vishal M Patel. Multi-level bottom-top and top-bottom feature fusion for crowd counting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1002–1012, 2019.

[21] 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*, pages 4036–4045, 2019.

[22] Qian Wang and Toby P Breckon. Segmentation guided attention network for crowd counting via curriculum learning. *arXiv preprint arXiv:1911.07990*, 2019.

[23] 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*, pages 8198–8207, 2019.

[24] Feng Xiong, Xingjian Shi, and Dit-Yan Yeung. Spatiotemporal modeling for crowd counting in videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5151–5159, 2017.

[25] 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 IEEE International Conference on Computer Vision*, pages 8362–8371, 2019.

[26] 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 IEEE International Conference on Computer Vision*, pages 8382–8390, 2019.

[27] 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 IEEE International Conference on Computer Vision*, pages 952–961, 2019.

[28] 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 IEEE International Conference on Computer Vision*, pages 6788–6797, 2019.

[29] 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*, pages 589–597, 2016.

[30] 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*, pages 12736–12745, 2019.

[31] Zhikang Zou, Huiliang Shao, Xiaoye Qu, Wei Wei, and Pan Zhou. Enhanced 3d convolutional networks for crowd counting. *arXiv preprint arXiv:1908.04121*, 2019.