

Hybrid attention network based on progressive embedding scale-context for crowd counting

Fusen Wang^{a,b}, Jun Sang^{a,b,*}, Zhongyuan Wu^{a,b}, Qi Liu^{a,b}, Nong Sang^c

^a Key Laboratory of Dependable Service Computing in Cyber Physical Society of Ministry of Education, Chongqing University, Chongqing 400044, China

^b School of Big Data & Software Engineering, Chongqing University, Chongqing 401331, China

^c School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan 430000, China

Abstract

The existing crowd counting methods usually adopted attention mechanism to tackle background noise, or applied multi-level features or multi-scales context fusion to tackle scale variation. However, these approaches deal with these two problems separately. In this paper, we propose an Hybrid Attention Network (HAN) by employing Progressive Embedding Scale-context (PES) information, which enables the network to simultaneously suppress noise and adapt head scale variation. We build the hybrid attention mechanism through paralleling spatial attention and channel attention module, which makes the network to focus more on the human head area and reduce the interference of background objects. Besides, we embed certain scale-context to the hybrid attention along the spatial and channel dimensions for alleviating these counting errors caused by the variation of perspective and head scale. Finally, we propose a progressive learning strategy through cascading multiple hybrid attention modules with embedding different scale-context, which can gradually integrate different scale-context information into the current feature map from global to local. Ablation experiments provides that the network architecture can gradually learn multi-scale features and suppress background noise. Extensive experiments demonstrate that HANet obtain state-of-the-art counting performance on four mainstream datasets.

Keywords: Crowd counting, Hybrid attention, Progressively embedding scale-context, Density map estimation

1. Introduction

Crowd counting based on deep learning has obtained widespread attention due to its significant applications in many large conferences, sports events, public transportation, etc. Its purpose is to count pedestrians in single image by training a density regression network under fully supervised learning. However, since

*Corresponding author

Email address: jsang@cqu.edu.cn (Jun Sang)

heavy scale variation and complex texture backgrounds exist in crowd images, the counting performance still remains room for improvement.

To tackle the heavy scale variation, many methods utilized multi-column networks with different convolution kernel sizes to extract multi-scales features for adapting to different head sizes [1, 2, 3, 4, 5]. In addition, some approaches [6], [7] encoded the context information of multiple receptive field sizes to handle the problems caused by the different density distribution and continuous changes of scale in the crowd image. On the other hand, to avoid the disturbance from background noise, some researchers employed visual attention mechanism to generate one attention weights map from the current or previous layer features through certain means for enhancing robustness to noise [5], [8], [9], [10], [11]. These methods simply divided the two problems of background noise and scale variation into different categories. Actually, the two problems are not independent while mutually affected and restricted each other, such as the network may not correctly distinguish small heads and leaves.

Some latest approaches also attempted to apply the attention mechanism to scale variety, but they neglect the obstruction of background noise. Besides, these works were often supported by some auxiliary tasks and needed to balance the calculation of multiple loss. This requested extra supervision data and the calculation of redundant loss [12, 13, 14, 15].

Above all, we propose one end-to-end hybrid attention network (HAN) of progressive embedding scale-context (PES) information to simultaneously suppress noise and adapt to head scale variation. This approach cascades several scale-context hybrid attention modules, which is composed of two parallel spatial attention and channel attention modules integrated with different head scale variation information from global to local.

To summarize, the main contributions of our work are outlined as follows:

- (1) We propose an end-to-end network by cascading multiple hybrid attention modules, which are composed of spatial attention and channel attention in parallel.
- (2) Context information of different scales is embedded in the proposed hybrid attention module through a progressive learning strategy, which can gradually adapt to people scale variety and suppress background noise.
- (3) The results demonstrate the proposed HANet achieves the state-of-the-art performance on the four benchmark datasets.

2. Related Work

In this section, we briefly review some mainstream related works on CNN-based crowd counting methods. These mainly include crowd counting approaches based on multi-column architecture and attention

mechanism.

2.1. Multi-column methods in crowd counting

Due to large perspective changes and varying resolution, the head size of pedestrians varies significantly in different positions of the crowd images. Many early successful works usually conducted multi-column network architecture to address this problem. For example, Zhang et al. [1] proposed Multi-Column Convolution Neural Network (MCNN) which was composed of three branches with different receptive field, aiming to tackle crowds with different densities. Following the above work, Sam et al. [3] trained a density classifier to adaptively select optimal regressor from three branch for different image patch (Switch CNN). Sindagi et al. [6] added two columns of pyramid branch on the basis of MCNN to extract global and local context information for accurate counting and generating high quality density map (CP-CNN). Boominathan et al. [2] employed two branches network with the same receptive filed of different depth, in which the shallow branch extracted small-scale features while deep branch was the opposite (CrowdNet). However, these methods brought a large amount of computation burden and were trained slowly to improve performance.

Therefore, some latest methods attempted new means to settle scale changes. Guo et al. [5] proposed multi-column dilated convolution network to capture different receptive fields required shifty head size without increasing redundant calculations. Liu et al. [7] proposed an end-to-end architecture that integrated multi-scale context information into current feature map in parallel (CAN).

2.2. Attention mechanism in crowd counting

Recently, attention mechanism has been widely incorporated to many tasks, such as object detection [16], image classification [17] and etc. It also works for crowd counting. The traditional attention-based crowd counting method usually generated a weight map by activation function from the current or previous layer features to mask background noise. For instance, Zhu et al. [8] proposed a dual path scale fusion network, where one path was employed to generate attention map for suppressing noise, while the other path multiplied the extracted features with the attention map to generate high quality density map (SFANet). Sindagi et al. [10] proposed a hierarchical attention-based network, which consisted of a spatial attention module and several global attention modules (HACNN).

Besides, some crowd counting methods of attention-based was aiming to adapt to the people head scale variety [12, 13, 14, 15]. Hossain et al. [12] proposed a scale-context attention network to focus on people density variation by adding two auxiliary branch namely Global scale attention and Local scale attention (SAAN). Jiang et al. [15] proposed multi-task density-aware network, which jointly trained density-level classification and density map estimation. The density-level classification utilized high-level semantic information to guide density map estimation (DensityCNN).

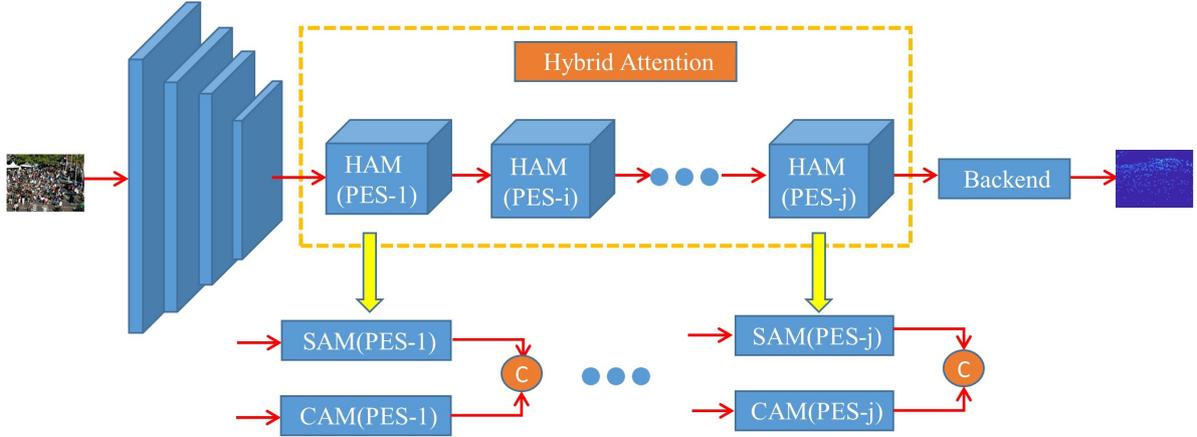


Figure 1: Overview of the proposed HANet architecture. It consists of Backbone (VGG16-BN first ten layers), Hybrid Attention (HA) that includes several cascade hybrid attention modules of progressively embedding scale-context (HAM-PES-K), and Backend.

3. Proposed Approaches

3.1. Overview of network architecture

To make the model better mask background noise on the basis of adapting to people scale variation, we design a new attention network of progressively embedding scale-context, namely HANet. In this section, we firstly illustrate the overview architecture of our proposed HANet and then introduce each component in detail. HANet includes three modules: (1) backbone; (2) hybrid attention; (3) backend, as shown in Figure 1.

3.2. Backbone module

VGG16 [18] has achieved excellent performance in crowd counting and other computer vision tasks such as object detection, classification, etc. As W-net [19] indicated that VGG-BN is considered as the relatively better baseline for crowd counting task, our backbone module adopts the first ten layers of the pretrained VGG16-BN from ImageNet as front-end network to extract rich features. Besides, since pooling layer results in losing vast location details and spatial information, we remove the final two pooling layers, which can ensure that the output of the network is $1/8$ of the original image resolution.

3.3. Hybrid attention

In this module, we introduce the proposed hybrid attention (HA), which includes several cascade hybrid attention modules (HAM) of progressively embedding scale-context (PES) information with different scales. As illustrated in Figure 1, the input of the first hybrid attention module (HAM) is the backbone’s output with the size of $C \times H \times W$, and the input of the subsequent HAM is the output of the previous HAM. We

design the HAM similar with DANet [20] through paralleling spatial attention (SAM) and channel attention (CAM) streams, except that our modules are designed with progressively embedding scale-context (PES) information in multi-scales rather than non-local mechanism. In other words, the multi-scales information is separately integrated into several hybrid attention modules along with the spatial and channel dimensions. The process can be formulated as follows:

$$x_i = \begin{cases} \mathcal{F}_{vgg}(X), & \text{for } i = 0 \\ \mathcal{F}[SAM_{PES(i)}^{\theta_i}, CAM_{PES(i)}^{\phi_i}](x_{i-1}), & \text{for } i > 0 \end{cases} \quad (1)$$

Where given an image X , x_0 denotes the Hybrid attention’s input from the output of the first ten layers of a pretrained VGG16-BN network; x_i ($i > 0$) denotes hybrid attention module of progressively embedding scale-context implemented by function \mathcal{F} , which represents the channel-wise concatenation operation for combining $SAM_{PES(i)}^{\theta_i}$ and $CAM_{PES(i)}^{\phi_i}$; the output x_i of the previous HAM is the input of the next HAM x_{i+1} ; x_n is the final result of entire Hybrid attention.

In CANet’s method [7], it simply concatenates context information of different scales with current features, which ignores that the scale variations are continuous and smooth, resulting in the network to fall into a partial solution of a certain scale. Different from CANet, our PES mechanism gradually embeds contextual information of different scales into several cascade hybrid attention modules from global to local. The HAM’s components SAM and CAM are depicted in detail in Figure 2.

3.3.1. Spatial attention module

Different from the previous scheme [21], [22], our spatial attention module (SAM) not only encodes the probability value of the head appearing in each position of the crowd image, but also incorporates certain scale-context information into each spatial attention module (SAM). The architecture of SAM(PES-K) is shown detailly at the top of Figure 2. The output of SAM is defined by the following equation:

$$SAM_{PES(K)} = G_3(x, \theta) * \sigma[G_1(G_3(x, \theta) + U(A_K(x)))] \quad (2)$$

where for input x with the size of $C \times H \times W$, it is fed into adaptive average pooling layer of scale K (A_K) to obtain $C \times K \times K$ scale-context blocks; U is the up-sample operation with bilinear interpolation to restore the scale-context blocks to the size of input x ; then add residual connection to the result of the Up step by G_3 operation, which represents a convolution layer of kernel size of 3×3 ; G_1 denotes a convolutional layer of kernel size of 1×1 for channel dimension reduction; Finally, a weight map generated by Sigmoid function (σ) from the previous output is multiplied with the residual features $G_3(x, \theta)$.

3.3.2. Channel attention module

The previous channel attention attempts to strengthen important channels (foreground) and suppress the unnecessary ones (background) [23]. However, the lack of scale-context information limits the feature

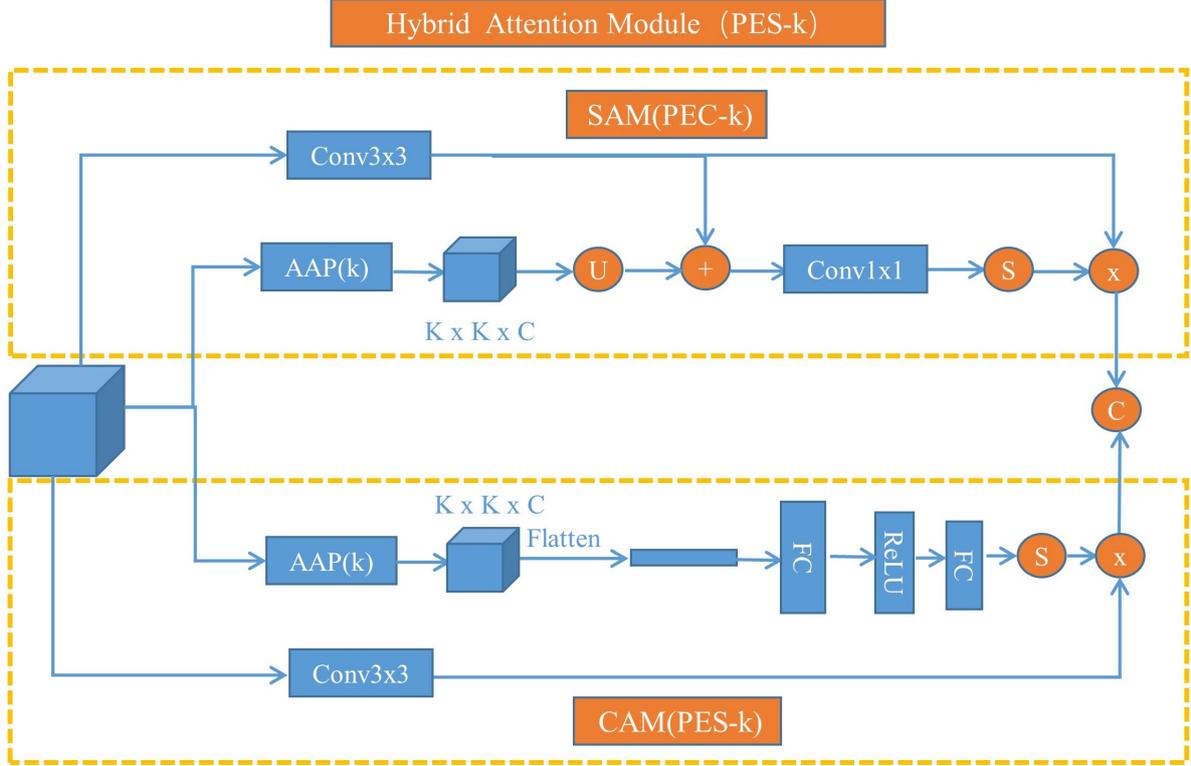


Figure 2: The detailed architecture of the Hybrid Attention Module of Progressively Embedding Scale-Context (PES-K). It includes two components, Spital Attention Module (PES-K) and Channel Attention Module (PES-K). AAP(k) indicates the adaptive average pooling of features into $K \times K \times C$ blocks; Conv $R \times R$, Flatten, and FC represent convolution with size of $R \times R$, flatten feature blocks of size $K \times K \times C$ into a $1 \times 1 \times K^2 C$ vector and fully connection layer, respectively; U, S, ReLU, C denotes severally bilinear interpolation, sigmoid and relu activation function, and channel-wise concatenation operation.

expression ability of the module. For example, it is difficult to distinguish small heads from leaves. We embed a scale-context blocks in CAM that is similar to SAM. The specific architecture of CAM(PES-K) is depicted at the bottom of Figure 2, and its mathematical expression formula is as follows:

$$CAM_{PES(K)} = G_3(x, \theta) * \sigma[F_C(F_L(A_K(x)))] \quad (3)$$

where the main operations are similar to the above SAM except F_C, F_L ; F_L expresses that the $C \times K \times K$ scale-context blocks from A_K stream are flattened into vectors with a size of $1 \times 1 \times K^2 C$; equivalent to [23], F_C consists of fully connection 1, ReLU activate function, and fully connection 2 three layers, which are employed to learn across-channels interactive and dependence related information. In the end, obtain a channel weight map by Sigmoid function multiplied with the residual features $G_3(x, \theta)$.

We conduct the hybrid attention module (HAM) by concatenation the above SAM and CAM of embedding K level scale-context information. Then, we cascade several HAM through a progressive learning

strategy as the overall hybrid attention architecture. The advantage of this strategy is to allow the model to gradually learn features of different scales instead of the previous multi-scale feature concatenating. Inspired by [7], we get the best results by adopting $K = 4$ different scales with 1, 2, 3, 6 in light of the trade-off between resource cost and performance. Furthermore, we adopt the embedding means of scale-context information from global to local, which is superior to the opposite approach. This will be validated in ablation studies of Sec. V.

3.4. Backend

The output of the above hybrid attention is fed to backend to generate the final density map through a set of convolution layers. These architectures are designed as follows: Conv (512,256,3)-BN-ReLU, Conv (256,128,3)-BN-ReLU, Conv (128,64,3)-BN-ReLU, Conv (64,1,1). During the training phase, we leverage the pixel-level Euclidean Loss (MSE) to measure the error between estimated density map and ground truth:

$$\mathcal{L}_{MSE} = \frac{1}{N} \sum_{k=1}^N \|D(X_k^{(i,j)}; \theta) - GT_k^{(i,j)}\|_2^2 \quad (4)$$

where N is the batch size of input’s images; $D(X_k^{(i,j)}; \theta)$ is the estimated density map at pixel location (i, j) for image X_k with parameters θ ; $GT_k^{(i,j)}$ is the corresponding ground truth for image X_k .

4. Implementation Details

In this section, we first introduce the four main-stream datasets as well as the method of generating ground truth. Sequentially, the data augmentation and training details are given respectively.

4.1. Ground Truth Generation

Similar to previous method [1], we blur the annotations point of crowd images by Gaussian kernel $G_{\mu,\sigma}$ with standard deviation μ and kernel size σ to estimate the size of people head. If there exist an annotated point at pixel x_i , it is denoted as $\delta(x - x_i)$. Then produce density maps through convolving $\delta(x - x_i)$. Its mathematical expression formula is as follows:

$$GT(x) = \sum_{i=1}^N \delta(x - x_i) \times G_{\mu,\sigma} \quad (5)$$

where N is the total number of annotated people in a crowd image. For the datasets with high congested scenes, we utilize geometry-adaptive kernel to generate density map while the fixed Gaussian kernel is used in several dataset of relatively sparse scenes. The method of producing the ground truth map in different datasets is illustrated in Table 1.

Table 1: Methods of generating ground truth map on different datasets.

Datasets	Method
ShanghaiTech Part A [1]	Fixed: $\mu = 15, \sigma = 4$
ShanghaiTech Part B [1]	Fixed: $\mu = 15, \sigma = 4$
UCF-QNRF [24]	Geometry-adaptive
UCF-CC-50 [25]	Fixed: $\mu = 15, \sigma = 4$

4.2. Datasets

4.2.1. ShanghaiTech

The ShanghaiTech [1] dataset is divided into two parts: Part A and Part B, which contains 1198 images with a total of 330,165 annotated heads. Part A includes 300 training images and 182 test images randomly downloaded from the Internet, where the resolutions of each image are greatly different. Part B is composed of 400 training images and 316 test images taken from streets on Shanghai, the image resolutions of which is 768×1024 .

4.2.2. UCF-CC-50

The UCF-CC-50 [25] is a very challenging crowd counting dataset due to its small size, in which the number of pedestrians each image ranges from 94 to 4,543 with an average number of 1,280 persons. It only includes 50 images of significantly congested scenes with a total of 63,974 head annotations. Since it is a small size dataset, we perform 5-fold cross validation by randomly selecting images to train and test our proposed approach.

4.2.3. UCF-QNRF

As the largest crowd dataset, the UCF-QNRF dataset [24] includes 1535 images of different scenarios from Internet with 1,251,642 annotations, among which it is divided into train and test set of 1201 and 334 images respectively. The number of pedestrians in each image varies from 49 to 12,865. Furthermore, the image resolutions are very large and its scale varied dramatically comparing other datasets. It is a challenging dataset for crowd counting due to the high-count images and a wider variety of scenes.

4.3. Data Augmentation

In training phase, we randomly crop M image patches with size of $m \times m$ pixels from the original image to ensure our network can be multi-batch trained and promote performance at a lower time cost. Table 2 details the configuration of M and m on different datasets. Also, we ensure that these M image patches are different from each other for avoiding the patch overlapping. In addition, for the diversity of training data,

Table 2: The crop size and number of image patches on different datasets.

Datasets	M	m
ShanghaiPart A [1]	4	128
ShanghaiPart B [1]	4	256
UCF-CC-50 [25]	8	128
UCF-QNRF [24]	8	128

we change the RGB image into gray with probability of 0.2, and randomly horizontal flip the image with probability 0.5. These means of data augmentation are employed in each iteration of the training process.

4.4. Training details

All experimental training and evaluation are implemented on the platform of PyTorch [26] with a NVIDIA Tesla k80 GPU. The baseline of the proposed HANet is leveraged from the first ten layers of pretrained VGG16-BN on ImageNet, and the other convolution layers are randomly initialized by Gaussian distributions with a mean 0 and standard deviation of 0.01. We use SGD optimizer with learning rate of 1e-4 and weight decay of 5e-4 to train our model by minimizing the loss function Eq. (4). The batch size is set to 8 or 16 according to GPU computing power and the number of iterations is set to 2000.

5. Experiments and Results

In this section, we give the evaluation metrics and compare the results of our method with state-of-the-art methods. In the end, we preform extensive ablation experiments to validate the effectiveness of our method.

5.1. Evaluation metrics

To evaluate the accuracy of our approach, the mean absolute error (MAE) and the mean squared error (MSE) are adopted as metrics. Specifically, equations are defined as:

$$MAE = \frac{1}{N} \sum_{i=1}^N |C_i^{ES} - C_i^{GT}| \quad (6)$$

$$MSE = \sqrt{\frac{1}{N} \sum_{i=1}^N |C_i^{ES} - C_i^{GT}|^2} \quad (7)$$

where N is the number of test images; C_i^{ES} is the estimated count of the i -th image, which can be calculated by integrating the estimated density map; C_i^{GT} is the corresponding ground truth map of the i -th image.

Table 3: Estimation errors on ShanghaiTech dataset.

Method	Part A		Part B	
	MAE	MSE	MAE	MSE
MCNN [1]	110.2	173.2	26.4	41.3
SANet [27]	67.0	104.5	8.4	13.6
SCAR [9]	66.3	114.1	9.5	15.2
SFCN [28]	64.8	107.5	7.6	13.0
TEDNet [29]	64.2	109.1	8.2	12.8
DADNet [5]	64.2	99.9	8.8	13.5
HACNN [10]	62.9	94.9	8.1	13.4
CAN [7]	62.3	100.0	7.8	12.2
MSANet [13]	62.1	98.5	7.6	12.4
MMNet [30]	60.8	99.0	7.6	11.7
RANet [31]	59.4	102.0	7.9	12.9
PaDNet [14]	59.2	98.1	8.1	12.2
PGCNet [32]	57.0	86.0	8.8	13.7
HANet (ours)	54.9	91.2	6.8	11.5

5.2. Comparisons with state-of-the-art

The extensive comparisons are conducted with state-of-the-art on four benchmark datasets to exhibit the effectiveness of our HANet. Furthermore, we also give the visual results.

ShanghaiTech. To evaluate the effectiveness of our HANet, we compare our model with 13 state-of-the-art approaches on the ShanghaiTech dataset. As shown in Table 3, the proposed HANet achieve the lowest MAE of 54.9 and a comparable MSE of 91.2 on Part A. For Part B, our model has also won two first places on MAE of 6.8 and MSE of 11.6.

UCF-QNRF. We further assess the counting performance of HANet through comparison with 6 state-of-the-art methods on the UCF-QNRF dataset, which has the highest crowd density images and contain a variety of perspective effects, density changes. The results are presented in Table 4. Our approach obtains the lowest MAE of 98 and the MSE of 179 closed to the best ones. This distinctly indicates that HANet has superior robustness against the density and scale variation.

UCF-CC-50. This is an extremely dense crowd dataset, which has the largest average number of each image. To evaluate the capability of HANet more precisely, we perform a 5-fold cross-validation that divides the dataset into 5parts, each including 40 training sets and 10 test sets. As shown in the Table 5, our proposed method obtains the lowest MAE with 195.2 and the second lowest MSE value of 268.6.

Table 4: Estimation errors on UCF-QNRF dataset.

Method	UCF-QNRF	
	MAE	MSE
MCNN [1]	277	426
HACNN [10]	118	180
TEDNet [29]	113	188
RANet [31]	111	190
CAN [7]	107	183
SFCN [28]	102	171
HANet (ours)	98	179

Table 5: Estimation errors on UCF-CC-50 DATASET.

Method	UCF-CC-50	
	MAE	MSE
MCNN [1]	377.6	509.1
SANet [27]	258.4	334.9
TEDNet [29]	249.4	354.5
CAN [7]	212.2	243.7
SFCN [28]	214.2	318.2
MSANet [13]	238.2	310.8
HANet (ours)	195.2	268.6

In the end, we take some samples from the test set of several crowd datasets for visual exhibition. As illustrated in Figure 3, the first, second and third column exhibit respectively the original images, ground truth and estimated density maps of HANet on ShanghaiTech Part A, B and UCF-QNRF dataset. From the figure, the results validate the performance of HANet in scenes with occluded backgrounds and severe scale variation.

5.3. Ablation Experiments

In this section, we conduct multigroup ablation experiments to evaluate the effectiveness of different components and strategies on ShanghaiTech A, B and UCF-CC-50 datasets.

5.3.1. Ablation analysis on model’s components

In this subsection, we demonstrate the performance of each module of our HANet on ShanghaiTech A dataset. Table 6 gives six different setting of network’s module combination to verify the effect.



Figure 3: An illustration of estimated density maps and crowd counts generated by proposed HANet. The first column shows three samples drawn from ShanghaiTech Part A, ShanghaiTech Part B and UCF-QNRF datasets. The second column shows the corresponding ground truth maps. The third column shows the density maps estimated by our HANet.

VGG16-BN first ten layers:the backbone of our model. We load the pretrained parameters on ImageNet to VGG16-BN. A 1x1 convolution kernel is immediately followed to generate the predicted density map;

VGG16-BN + Backend: The backend includes four convolution layers (Conv512-256-3, Conv256-128-3, Conv128-64-3, Conv64-1-1), which are added to the end of backbone to regress predicted density maps more accurately;

VGG16-BN + HAM(PES-I) + ... + HAM(PES-J) + Backend(HANet): Progressive add every hybrid attention module of embedding certain scale-context information to the end of backbone from global to local. Then Backend is employed to generate the estimated density map.

From the table, considering the number of parameters and calculation, the performance of our model gradually improves to 54.9/91.2 of MAE/MSE through continuously adding HAM(PES-K) to HAM(PES-1,2,3,6). This demonstrates the effectiveness of our approach.

In Figure 4, we visualize the estimated density map generated by different components of the proposed model in ShanghaiTech A dataset. The first row represents the original image and the corresponding ground

Table 6: Estimation errors for different components of the proposed method on ShanghaiTech A dataset.

Module	MAE	MSE	Parameters(M)	Flops(G)
VGG16-BN	68.3	108.7	—	—
VGG16-BN+Backend	65.4	102.2	—	—
VGG16-BN + HAM(PES-1) + Backend	61.0	97.4	11.78	5.57
VGG16-BN + HAM(PES-1,2) + Backend	57.9	91.5	14.56	6.17
VGG16-BN + HAM(PES-1,2,3) + Backend	56.6	95.1	17.68	6.78
VGG16-BN + HAM(PES-1,2,3,6) + Backend	54.9	91.2	22.56	7.45

Table 7: Estimation errors of the fusion means on ShanghaiTech A dataset.

Module	MAE	MSE
VGG16-BN + HAM(PES-3,2,1) + Backend	57.8	95.3
VGG16-BN + HAM(PES-1,2,3) + Backend	56.6	95.1
VGG16-BN + HAM(PES-6,3,2,1) + Backend	56.6	93.2
VGG16-BN + HAM(PES-1,2,3,6) + Backend	54.9	91.2

truth, respectively. The second-third rows separately denote the estimated density map produced by multiple hybrid attention of progressively embedding scale-context information, i.e., EST Count (1,2,3,6) represents the number of people estimated by the scale combination as 1,2,3,6 in hybrid attention of progressively embedding scale-context, and the rest is similar. From the figure, we find that the counting performance has been continuously improved close to ground truth through gradually adding attention module of embedding scale-context.

5.3.2. Ablation analysis on HAM(PES)’s fusion methods

Table 7 displays the comparison of two fusion methods for hybrid attention module of progressively embedding scale-context on ShanghaiTech A dataset, which represent the approaches of local to global (i.e., PES 3,2,1 or PES 6,3,2,1) and global to local (i.e., PES 1,2,3 or PES 1,2,3,6). We find that the model prefers to greatly the scale-learning means of the latter. From the results in the first and second rows in the Table VII, the fusion approach of PES 1,2,3 achieves better performance of MAE 56.6 and MSE 95.1. Similarly, in the third and fourth rows in the Table 7, the MAE and MSE of fusion method of PES 1,2,3,6 can be reduced from 56.6 and 93.2 to 54.9 and 91.2.

5.3.3. Ablation analysis on size of image patches.

Due to the progressively embedding scale-context (PES) of hybrid attention module, the size of the image patch also seriously affects the counting performance of the model. For datasets with different population

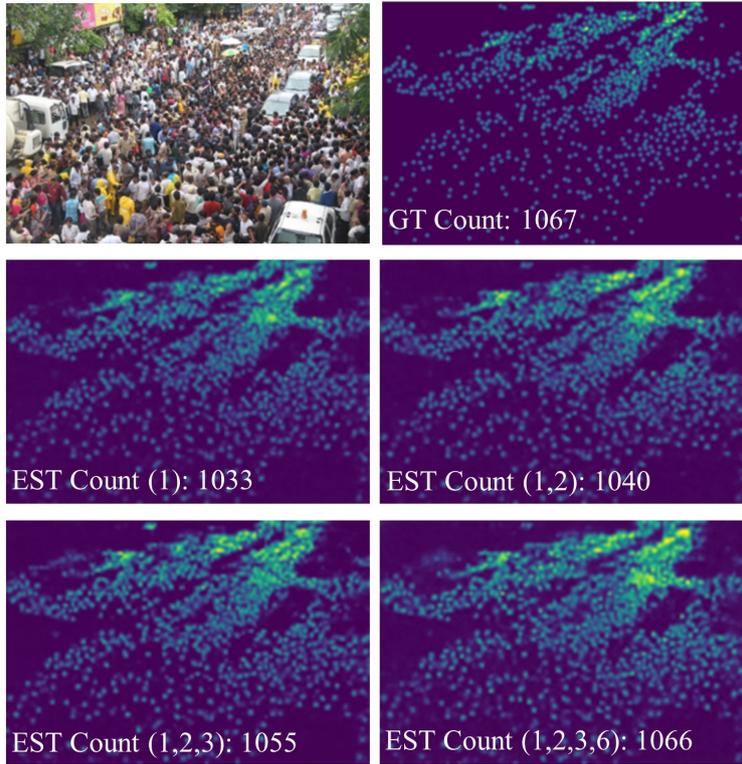


Figure 4: Estimated density maps generated different components of the model in ShanghaiTech A dataset for comparison the effectiveness of each module.

area densities, the size of image patches that can improve the best performance of the network are also different. As shown in Table 8, we compare the impact of different image patch’s sizes for model performance on ShanghaiTech B and UCF-CC-50 datasets. This also demonstrates the experimental configuration in Table 2.

From the above table, we observed that the size of image patch has a significant influence for experimental results on the datasets of different population density distribution. Due to the huger number of heads and the wider crowd density on UCF-CC-50 dataset, which represent the smaller image patches usually have

Table 8: The effects of different patch sizes with HAM(PES) module

Size	Part B		UCF-CC-50	
	MAE	MSE	MAE	MSE
128x128	7.2	13.6	195.2	268.6
192x192	7.0	12.5	211.8	302.1
256x256	6.8	11.5	201.8	277.6

contained more rich heads feature information with a similar scale. Therefore, the patch of 128x128 (MAE of 195.2 and MSE of 268.6) achieved better performance than the larger patch of 192x 192 (MAE of 211.8 and MSE of 302.1) and 256x256 (MAE of 201.8 and MSE of 277.6). However, since ShanghaiTech B includes a relatively small population and a lower population density, this requires larger image patches to capture more available people scale-context information to improve the effect of model. In our experiments, the patches of 256x256 obtain lower MAE of 6.8 and MSE of 11.5 than patches of 128x128 (MAE of 7.2 and MSE of 13.6) and 192x192 (MAE of 7.0 and MSE of 12.5).

6. Conclusion

In this paper, we propose a hybrid attention network based on progressive embedded scale-context information for crowd counting. To distinguish human head and background more accurately, we gradually embed certain scale-context information to attention mechanism, which enables the network suppress noise and adapt to head scale changes. Our method is validated on four mainstream datasets, and superior counting performance is obtained in comparison with current state-of-the-art approaches.

Acknowledgment

This work was supported by National Natural Science Foundation of China (No. 61971073).

References

- [1] Y. Zhang, D. Zhou, S. Chen, S. Gao, Y. Ma, Single-image crowd counting via multi-column convolutional neural network, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 589–597. doi:10.1109/CVPR.2016.70.
- [2] L. Boominathan, S. S. Kruthiventi, R. V. Babu, Crowdnet: A deep convolutional network for dense crowd counting, in: Proceedings of the 24th ACM international conference on Multimedia, 2016, pp. 640–644.
- [3] D. B. Sam, S. Surya, R. V. Babu, Switching convolutional neural network for crowd counting, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 4031–4039. doi:10.1109/CVPR.2017.429.
- [4] J. Liu, C. Gao, D. Meng, A. G. Hauptmann, Decidenet: Counting varying density crowds through attention guided detection and density estimation, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 5197–5206. doi:10.1109/CVPR.2018.00545.
- [5] D. Guo, K. Li, Z.-J. Zha, M. Wang, Dadnet: Dilated-attention-deformable convnet for crowd counting, in: Proceedings of the 27th ACM International Conference on Multimedia, 2019, pp. 1823–1832.
- [6] V. A. Sindagi, V. M. Patel, Generating high-quality crowd density maps using contextual pyramid cnns, in: 2017 IEEE International Conference on Computer Vision (ICCV), 2017, pp. 1879–1888. doi:10.1109/ICCV.2017.206.
- [7] W. Liu, M. Salzmann, P. Fua, Context-aware crowd counting, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 5094–5103. doi:10.1109/CVPR.2019.00524.
- [8] L. Zhu, Z. Zhao, C. Lu, Y. Lin, Y. Peng, Dual path multi-scale fusion networks with attention for crowd counting, arXiv preprint arXiv:1902.01115.

- [9] J. Gao, Q. Wang, Y. Yuan, Scar: Spatial-/channel-wise attention regression networks for crowd counting, *Neurocomputing* 363 (2019) 1–8.
- [10] V. A. Sindagi, V. M. Patel, Ha-ccn: Hierarchical attention-based crowd counting network, *IEEE Transactions on Image Processing* 29 (2020) 323–335. doi:10.1109/TIP.2019.2928634.
- [11] N. Liu, Y. Long, C. Zou, Q. Niu, L. Pan, H. Wu, Adcrowdnet: An attention-injective deformable convolutional network for crowd understanding, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 3220–3229. doi:10.1109/CVPR.2019.00334.
- [12] M. Hossain, M. Hosseinzadeh, O. Chanda, Y. Wang, Crowd counting using scale-aware attention networks, in: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), 2019, pp. 1280–1288. doi:10.1109/WACV.2019.00141.
- [13] R. R. Varior, B. Shuai, J. Tighe, D. Modolo, Multi-scale attention network for crowd counting, arXiv preprint arXiv:1901.06026.
- [14] Y. Tian, Y. Lei, J. Zhang, J. Z. Wang, Padnet: Pan-density crowd counting, *IEEE Transactions on Image Processing* 29 (2020) 2714–2727. doi:10.1109/TIP.2019.2952083.
- [15] X. Jiang, L. Zhang, T. Zhang, P. Lv, B. Zhou, Y. Pang, M. Xu, C. Xu, Density-aware multi-task learning for crowd counting, *IEEE Transactions on Multimedia* 23 (2021) 443–453. doi:10.1109/TMM.2020.2980945.
- [16] D. Yoo, S. Park, J. Lee, A. S. Paek, I. S. Kweon, Attentionnet: Aggregating weak directions for accurate object detection, in: 2015 IEEE International Conference on Computer Vision (ICCV), 2015, pp. 2659–2667. doi:10.1109/ICCV.2015.305.
- [17] C. C. et al, Look and think twice: Capturing top-down visual attention with feedback convolutional neural networks, in: Proceedings of the IEEE international conference on computer vision, 2015, pp. 2956–2964.
- [18] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- [19] V. K. Valloli, K. Mehta, W-net: Reinforced u-net for density map estimation, arXiv preprint arXiv:1903.11249.
- [20] J. Fu, J. Liu, H. Tian, Y. Li, Y. Bao, Z. Fang, H. Lu, Dual attention network for scene segmentation, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 3141–3149. doi:10.1109/CVPR.2019.00326.
- [21] S. Woo, J. Park, J.-Y. Lee, I. S. Kweon, Cbam: Convolutional block attention module, in: Proceedings of the European conference on computer vision (ECCV), 2018, pp. 3–19.
- [22] Z. Wang, J. Xu, L. Liu, F. Zhu, L. Shao, Ranet: Ranking attention network for fast video object segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 3978–3987.
- [23] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, in: 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2018, pp. 7132–7141. doi:10.1109/CVPR.2018.00745.
- [24] H. I. et al, Composition loss for counting, density map estimation and localization in dense crowds, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 532–546.
- [25] H. Idrees, I. Saleemi, C. Seibert, M. Shah, Multi-source multi-scale counting in extremely dense crowd images, in: 2013 IEEE Conference on Computer Vision and Pattern Recognition, 2013, pp. 2547–2554. doi:10.1109/CVPR.2013.329.
- [26] A. P. et al, Automatic differentiation in pytorch.
- [27] X. Cao, Z. Wang, Y. Zhao, F. Su, Scale aggregation network for accurate and efficient crowd counting, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 734–750.
- [28] Q. Wang, J. Gao, W. Lin, Y. Yuan, Learning from synthetic data for crowd counting in the wild, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 8190–8199. doi:10.1109/CVPR.2019.00839.
- [29] X. Jiang, Z. Xiao, B. Zhang, X. Zhen, X. Cao, D. Doermann, L. Shao, Crowd counting and density estimation by trellis encoder-decoder networks, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 6126–6135. doi:10.1109/CVPR.2019.00629.
- [30] L. Dong, H. Zhang, Y. Ji, Y. Ding, Crowd counting by using multi-level density-based spatial information: A multi-scale

cnn framework, *Information Sciences* 528 (2020) 79–91.

- [31] A. Zhang, J. Shen, Z. Xiao, F. Zhu, X. Zhen, X. Cao, L. Shao, Relational attention network for crowd counting, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 6787–6796. doi:10.1109/ICCV.2019.00689.
- [32] Z. Yan, Y. Yuan, W. Zuo, X. Tan, Y. Wang, S. Wen, E. Ding, Perspective-guided convolution networks for crowd counting, in: 2019 IEEE/CVF International Conference on Computer Vision (ICCV), 2019, pp. 952–961. doi:10.1109/ICCV.2019.00104.