

## Spatial-related and Scale-aware Network for Crowd Counting

Lei Li<sup>1</sup>, Yuan Dong<sup>1</sup>, Hongliang Bai<sup>2</sup> <sup>1</sup>Beijing University of Posts and Telecommunications <sup>2</sup>Beijing Faceall Co

## I. Introduction

· Crowd counting aims to estimate the number of people in images. The challenge of crowd counting task lies in the variations of scale, different perspective, cluttering, background noise and occlusions. In this paper, we propose a learnable spatial attention module which can get the spatial relations to diminish the negative impact of backgrounds. Besides, a dense hybrid dilated convolution module is also brought up to preserve information derived from varied scales. With these two modules, our network can deal with the problems caused by scale variance and background interference. Our method can achieve state-of-the-art results on three representative crowd counting benchmarks.



• Overview of the proposed framework: We use a standard image classification network VGG-19 as our backbone, with the last pooling and the subsequent fully connected layers removed. Our LSAM and DHDCM are added after the backbone. Next, we upsample the output after the above modules by two times with bilinear interpolation, and then feed it to a regression head consisting of three convolutional layers to get the density map.



- *LSAM:* To get the spatial relations in the feature maps, we use a learnable convolution operation to capture the spatial attentions. The lower branch consists of three convolutional layers and a softmax layer. The convolutional layers are designed to get the spatial attentions.
- *DHDCM:* To solve the problem of scale variation, we use dilated convolutional layers to get the most representative features of people across different scales. A module to densely connect a set of dilated convolutional layers, so that the generated multi-scale features can cover a denser scale range without significantly increasing the parameters.

## III. Expriments

• *Results:* The largely varied heads of different sizes can be captured, background and crowd region can be easily distinguished. Apart from that, the estimated numbers of people is close to the ground truth, proving the robustness and accuracy of our model, not only in congested crowd scenes but also sparse scenes.



- *Ablation Study of LSAM and DHDCM:* Simultaneously take the background and scale influence into consideration can respectively bring improvements compared with the baseline method.
- *Comparison with SOTA methods:* All the results in the three datasets prove that the proposed LSAM and DHDCM are effective to optimize the performance, not only in dense scenes but also in sparse scenes.

ABLATION STUDY OF LSAM AND DHDCM ON UCF-QNRF BENCHMARK.					COMPARISON WITH THE STATE OF THE ART METHODS ON UCF-QNRF BENCHMARK.			
LSAM D	HDCM	MA	E M	SE	Method	MAE	MSE	E
- - -	- - ~	98. 96. 95.	6 17 2 16 4 16	2.8 8.5 4.3	MCNN[1] SwitchCNN[6] CL[2] S-DCNet[28]	277 228 132 104.4	426 445 191 176.	1
	~	93.	2 15	8.2	Ours	93.2	158.	2
Comparison with the state of the Art methods on ShanghaiTech benchmark.					Comparison with the state of the Art methods on UCF_CC_50 benchmark.			
Method	Par	PartA		rtB	Method	M	AE	MSE
	MAE	MSE	MAE	MSE	MCNN[1]	37	7.6	509.1
MCNN[1] SwitchCNN[6]	90.4	173.2	26.4	41.3	SwitchCNN[6]	31	8.1	439.2
IC-CNN[27]	68.5	116.2	10.7	16.0	CSRNet[9]	26	6.1	397.5
CSRNet[9]	68.2	115.0	10.6	16.0	SANet[7]	25	8.4	334.9
BL[26]	64.5	104.0	7.9	13.3	ADCrowdNet[1]	21 25	7.1	363.5
BL+[26]	62.8	101.8	7.7	12.7	BI +[26]	22	93	308.2
ADCrowdNet[12] Ours	63.2 60.95	98.9 <b>97.55</b>	7.7 7.5	12.9 12.6	Ours	22	0.5	302.9

## IV. Conclusions

- *a.* Learnable Spatial Attention Module can get spatial attentions.
- **b.** Dense Hybrid Dilated Convolutional Module can solve the scale variation problem.
- *c*. Both modules can be transferred to other related task.
- d. The state-of-the-art results on all three datasets.