



# CSRNet: Dilated Convolutional Neural Networks for Understanding the Highly Congested Scenes

Yuhong Li<sup>1,2</sup>, Xiaofan Zhang<sup>1</sup>, Deming Chen<sup>1</sup>

<sup>1</sup>ECE, University of Illinois at Urbana-Champaign, <sup>2</sup>Beijing University of Posts and Telecommunications



For full text ↑

## Motivations

### Major tasks for understanding the highly congested scenes

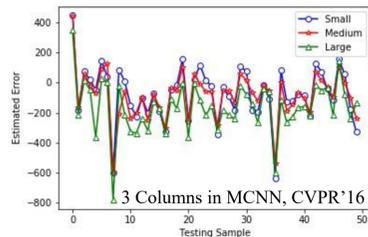
1. Count the crowd → get the number
2. Generate density map → get the crowd distribution
3. Understand the potential risks and make decisions
  - Total crowd number, Density patterns, Change rate, etc.



### Drawbacks of the state-of-the-art

- Use multi-column (MCNN) structures (branches/columns for different levels of density)
- Different columns are likely to learn similar features

E.g., MCNN, CVPR'16  
Crowdnet, ACM MM'16  
SwitchingCNN, CVPR'17



Method	Parameters	MAE	MSE
Col. 1 of MCNN	57.75k	141.2	206.8
Col. 2 of MCNN	45.99k	160.5	239.0
Col. 3 of MCNN	25.14k	153.7	230.2
MCNN Total	127.68k	110.2	185.9
A deeper CNN	83.84k	<b>93.0</b>	<b>142.2</b>

A deeper CNN (w/o branches) performs even better with less parameters

- MCNN-like structures require more efforts to train
- Hard to deploy on IoT devices (larger size with more parameters)

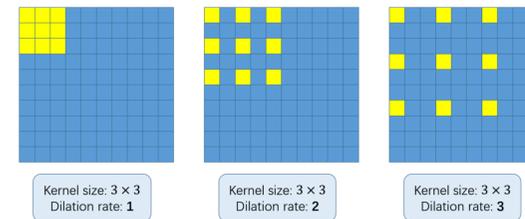
## Acknowledgement

This work was supported by the IBM-Illinois Center for Cognitive Computing System Research (C<sup>3</sup>SR) - a research collaboration as part of the IBM AI Horizons Network.

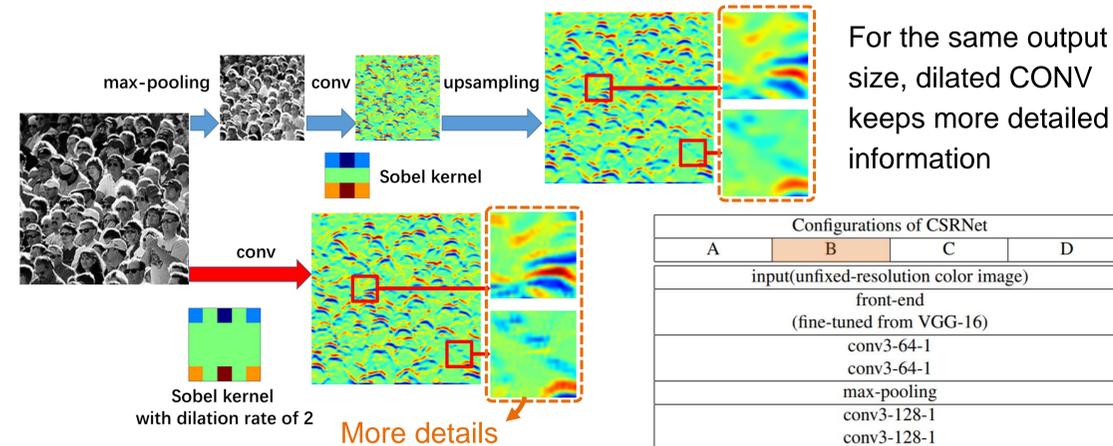
## Solution: CSRNet

### An end-to-end single-column structure delivering the best performance

- Front-end: CNN → Extract 2D features
- Back-end: Dilated CNN → Larger reception fields without losing accuracy (compared to pooling)



### Advantage of using dilated CONV instead of pooling



For the same output size, dilated CONV keeps more detailed information

### CSRNet Network configuration

conv(kernel size)-(# of filters)-(dilation rate)  
max-pooling layers are conducted over a 2x2 pixel window with stride 2.

Architecture	MAE	MSE
CSRNet A	69.7	116.0
CSRNet B	<b>68.2</b>	<b>115.0</b>
CSRNet C	71.91	120.58
CSRNet D	75.81	120.82

Evaluation in ShanghaiTech Part A

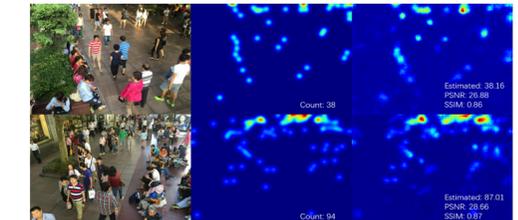
Configurations of CSRNet			
A	B	C	D
input(unfixed-resolution color image)			
front-end (fine-tuned from VGG-16)			
conv3-64-1			
conv3-64-1			
max-pooling			
conv3-128-1			
conv3-128-1			
max-pooling			
conv3-256-1			
conv3-256-1			
conv3-256-1			
max-pooling			
conv3-512-1			
conv3-512-1			
conv3-512-1			
back-end (four different configurations)			
conv3-512-1	conv3-512-2	conv3-512-2	conv3-512-4
conv3-512-1	conv3-512-2	conv3-512-2	conv3-512-4
conv3-512-1	conv3-512-2	conv3-512-2	conv3-512-4
conv3-256-1	conv3-256-2	conv3-256-4	conv3-256-4
conv3-128-1	conv3-128-2	conv3-128-4	conv3-128-4
conv3-64-1	conv3-64-2	conv3-64-4	conv3-64-4
conv1-1-1			

## Results

We demonstrate CSRNet on ShanghaiTech, UCF CC 50, WorldExpo'10, UCSD datasets for crowd counting, and TRANCOS dataset for car counting.

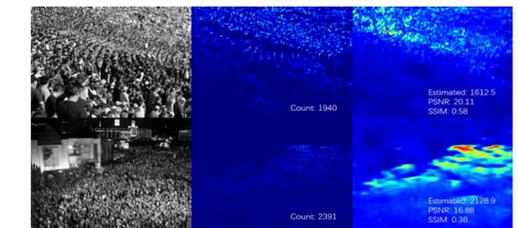
### Results in ShanghaiTech dataset

Method	Part_A		Part_B	
	MAE	MSE	MAE	MSE
Zhang <i>et al.</i> [3]	181.8	277.7	32.0	49.8
Marsden <i>et al.</i> [38]	126.5	173.5	23.8	33.1
MCNN [18]	110.2	173.2	26.4	41.3
Cascaded-MTL [39]	101.3	152.4	20.0	31.1
Switching-CNN [4]	90.4	135.0	21.6	33.4
CP-CNN [5]	73.6	<b>106.4</b>	20.1	30.1
CSRNet (ours)	<b>68.2</b>	115.0	<b>10.6</b>	<b>16.0</b>



### Results in UCF CC 50 dataset

Method	MAE	MSE
Idrees <i>et al.</i> [22]	419.5	541.6
Zhang <i>et al.</i> [3]	467.0	498.5
MCNN [18]	377.6	509.1
Marsden <i>et al.</i> [38]	338.6	424.5
Cascaded-MTL [39]	322.8	397.9
Switching-CNN [4]	318.1	439.2
CP-CNN [5]	295.8	<b>320.9</b>
CSRNet (ours)	<b>266.1</b>	397.5



### Results in WorldExpo'10 dataset

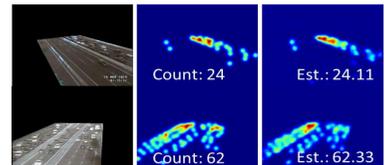
Method	Scce.1	Scce.2	Scce.3	Scce.4	Scce.5	Avg.
Chen <i>et al.</i> [46]	2.1	55.9	9.6	11.3	3.4	16.5
Zhang <i>et al.</i> [3]	9.8	14.1	14.3	22.2	3.7	12.9
MCNN [18]	3.4	20.6	12.9	13.0	8.1	11.6
Shang <i>et al.</i> [37]	7.8	15.4	14.9	11.8	5.8	11.7
Switching-CNN [4]	4.4	15.7	10.0	11.0	5.9	9.4
CP-CNN [5]	2.9	14.7	10.5	<b>10.4</b>	5.8	8.86
CSRNet (ours)	<b>2.9</b>	<b>11.5</b>	<b>8.6</b>	16.6	<b>3.4</b>	<b>8.6</b>

### Results in UCSD dataset

Method	MAE	MSE
Zhang <i>et al.</i> [3]	1.60	3.31
CCNN [20] CCNN	1.51	-
Switching-CNN [4]	1.62	2.10
FCN-rLSTM [24]	1.54	3.02
CSRNet (ours)	1.16	1.47
MCNN [18]	<b>1.07</b>	<b>1.35</b>

### Results in TRANCOS dataset

Method	GAME 0	GAME 1	GAME 2	GAME 3
Fiaschi <i>et al.</i> [47]	17.77	20.14	23.65	25.99
Lempitsky <i>et al.</i> [32]	13.76	16.72	20.72	24.36
Hydra-3s [20]	10.99	13.75	16.69	19.32
FCN-HA [24]	4.21	-	-	-
CSRNet (Ours)	<b>3.56</b>	<b>5.49</b>	<b>8.57</b>	<b>15.04</b>



Source code: <https://github.com/leeyeehoo/CSRNet>

## Conclusions

- The single-column structure is more efficient than MCNNs  
Less complexity & parameters, easier to train & reproduction
- We adapt dilated CONV to crowd counting applications  
Get larger receptive field without losing details
- We deliver the state-of-the-art performance in five datasets