



HANet: Hybrid Attention-aware Network for Crowd Counting

Xinxing Su^{*1}, Yuchen Yuan^{*2}, Xiangbo Su², Zhikang Zou², Shilei Wen², and Pan Zhou¹

¹Huazhong University of Science and Technology ²Baidu Inc.











Motivation

◆ Task

- Estimating the number of people
- Generating high-quality density maps



- non-uniform crowd distribution, noisy background, occlusions
- restricted to the 2D position-wise (or firstorder) attention
- independence between different supervisions



(a)

(a) Raw image; (b) Ground truth density map

(b)



Non-uniform crowd distribution, noisy background, occlusions

Overview



Our proposed HANet.

The green dotted line is the adaptive compensation loss (ACLoss).

Overview

Inspired by Kernel Pooling[26], we introduce the High-order Attention Module (HAM), which captures higher-ordered feature of the target via a kernel function t(x):

$$t(\mathbf{x}) = \langle W, \phi(\mathbf{x}) \rangle$$

= $\sum_{r=1}^{R} \langle W^r, \bigotimes_r \mathbf{x} \rangle = \langle W^1, \mathbf{x} \rangle + \sum_{r=2}^{R} \langle W^r, \bigotimes_r \mathbf{x} \rangle$

where $\langle \cdot, \cdot \rangle$ calculates the inner product, *R* is the maximum rank of the polynomial, W^r is a r-ranked learnable tensor.

Due to the sparsity of the CNN, in order to minimum the storage cost and risk of overfitting, we decompose W^r by [55]:

$$W^r = \sum_{d=1}^{D^r} \beta^{r,d} v_1^{r,d} \otimes v_2^{r,d} \otimes \cdots \otimes v_r^{r,d},$$

where W^r is decomposed into the summation of D^r tensors; v is the obtained tensor, and $\beta^{r,d}$ is the corresponding weight.

High-order Attention Module (HAM)

 $\beta^r =$



HAM with R=3

$$\begin{split} t(\mathbf{x}) \text{ is thus rewritten as:} \\ t(\mathbf{x}) &= \langle W^1, \mathbf{x} \rangle + \sum_{r=2}^R \left\{ \sum_{d=1}^{D^r} \beta^{r,d} v_1^{r,d} \otimes v_2^{r,d} \otimes \cdots \otimes v_r^{r,d}, \otimes_r \mathbf{x} \right\} \\ &= \langle W^1, \mathbf{x} \rangle + \sum_{r=2}^R \sum_{d=1}^{D^r} \beta^{r,d} \prod_{s=1}^r \langle v_s^{r,d}, \mathbf{x} \rangle \\ &= \langle W^1, \mathbf{x} \rangle + \sum_{r=2}^R \langle \beta^r, z^r \rangle \\ &\text{where } z^r = [z^{r,1}, z^{r,2}, \dots, z^{r,D^r}]^T \, \text{fll} z^{r,d} = \prod_{s=1}^r \langle v_s^{r,d}, \mathbf{x} \rangle, \\ \beta^r = [\beta^{r,1}, \dots, \beta^{r,D^r}]^T. \text{ Let } W^1 = v^1 \beta^1, \text{ where } v^1 \in \mathbb{R}^{C \times D^1}, \ \beta^1 \in \mathbb{R}^{1-q} \end{split}$$

 $\mathbb{R}^{D^1 \times C}$, t(x) is then transformed as: $t(\mathbf{x}) = \langle \beta^1, (v^1)^T \mathbf{x} \rangle + \sum_{r=2}^{R} \langle \beta^r, z^r \rangle = \sum_{r=1}^{R} \langle \hat{\beta}^r, \hat{z}^r \rangle.$

The Adaptation Compensation Loss (ACLoss)

Our supervision is:

$$\mathcal{L}_{total} = \sum_{l=1}^{B} \sum_{i=0}^{4} \| \hat{Y}_{l,i} - Y_{l,i} \|_{2}^{2},$$

where *B* is the batchsize, *i* is the index of the decoder, $\hat{Y}_{l,i}$ is the density map of image *l* decoded from module *i*, and $Y_{l,i}$ is the corresponding ground truth. Furthermore, the ACLoss can be updated as:

$$\begin{split} \mathcal{L}_{total} &= \sum_{l=1}^{B} \left(\sum_{i=1}^{3} \left(1 + w_{l,i} \right) \parallel \hat{Y}_{l,i} - Y_{l,i} \parallel_{2}^{2} \right. \\ &+ \parallel \hat{Y}_{l,0} - Y_{l,0} \parallel_{2}^{2} + \parallel \hat{Y}_{l,4} - Y_{l,4} \parallel_{2}^{2}), \end{split}$$

Where $w_{l,i} = |\varphi(\hat{Y}_{l,i-1}) - \varphi(Y_{l,i-1})|, \varphi(\cdot)$ is the Sigmoid function.

Experimental Results

Quantitative Results

EXPERIMENTAL RESULTS ON THE SHANGHAI PART A, SHANGHAI PART B AND UCF_CC_50 DATASETS

EXPERIMENTAL RESULTS ON THE UCF-QNRF DATASET

Method	Shangh	ai Part A	Shangh	ai Part B	UCF_CC_50		
wiethod	MAE↓	RMSE↓	MAE↓	RMSE↓	MAE↓	RMSE↓	
MCNN [17]	110.2	173.2	26.4	41.3	377.6	509.1	
CSRNet [6]	68.2	115.0	10.6	16.0	266.1	397.5	
ADCrowdNet [10]	66.1	102.1	7.6	13.9	257.9	357.7	
TEDnet [15]	64.2	109.1	8.2	12.8	249.4	354.4	
DADNet [12]	64.2	99.9	8.8	13.5	285.5	389.7	
CAN [11]	62.3	100.0	7.8	12.2	212.2	243.7	
ADMG [37]	64.7	97.1	8.1	13.6	239.8	319.4	
ANF [31]	63.9	99.4	8.3	13.2	250.2	340.0	
MBTTBF-SCFB [28]	60.2	94.1	8.0	15.5	233.1	300.9	
RANet [30]	59.4	102.0	7.9	12.9	239.8	319.4	
HANet (ours)	58.5	92.4	7.0	10.3	225.9	269.1	

Method	MAE↓	RMSE↓	
MCNN [17]	277.0	426.0	
Idress et al. [18]	132.0	191.0	
DADNet [12]	113.2	189.4	
TEDnet [15]	113.0	188.0	
RANet [30]	111.0	190.0	
CAN [11]	107.0	183.0	
ANF [31]	110.0	174.0	
ADMG [37]	101.0	176.0	
MBTTBF-SCFB [28]	97.5	165.2	
HANet (ours)	95.3	160.3	

Qualitative Results

Experimental Results



Top to bottom: input, ground truth density map, inference result of CSRNet[6] and that of HANet.

Experimental Results

Ablation Study

PERFORMANCE COMPARISONS ON SHANGHAITECH PART A WITH RESPECT TO THE ORDER AND THE PLACEMENT OF HAM.

Metric enco	encoder	E ₂			E ₃			E4					
	cheoder	R=1	R=2	R=4	R=6	R=1	R=2	R=4	R=6	R=1	R=2	R= 4	R=6
MAE	64.8	63.5	63.7	64.1	63.7	63.2	62.9	61.9	62.7	63.7	63.2	62.3	63.5

HAM And ACLoss:

base net	HAM	ACLoss	Shangh	ai Part A	UCF-QNRF		
			MAE↓	RMSE↓	MAE↓	RMSE↓	
\checkmark	-	-	62.6	104.9	100.1	180.2	
\checkmark	\checkmark	-	59.6	100.1	97.3	175.0	
\checkmark	-	\checkmark	61.6	100.9	96.4	171.1	
\checkmark	\checkmark	\checkmark	58.5	92.4	95.3	160.3	

Summary

Contributions

- Establishing a 3D attention mechanism that captures high-order statistics, both position-wisely and channel-wisely
- Exploiting the difference between the prediction and the ground truth from a higher-level supervision path as an additional attention guidance

THANK YOU FOR WATCHING