

CROSS-HEAD SUPERVISION FOR CROWD COUNTING WITH NOISY ANNOTATIONS

Mingliang Dai¹ Zhizhong Huang¹ Jiaqi Gao¹ Hongming Shan² Junping Zhang^{1†}

¹ Shanghai Key Lab of Intelligent Information Processing, School of Computer Science

² Institute of Science and Technology for Brain-inspired Intelligence

Fudan University, Shanghai 200433, China

ABSTRACT

Noisy annotations such as missing annotations and location shifts often exist in crowd counting datasets due to multi-scale head sizes, high occlusion, etc. These noisy annotations severely affect the model training, especially for density map-based methods. To alleviate the negative impact of noisy annotations, we propose a novel crowd counting model with one convolution head and one transformer head, in which these two heads can supervise each other in noisy areas, called **Cross-Head Supervision**. The resultant model, CHS-Net, can synergize different types of inductive biases for better counting. In addition, we develop a progressive cross-head supervision learning strategy to stabilize the training process and provide more reliable supervision. Extensive experimental results on ShanghaiTech and QNRF datasets demonstrate superior performance over state-of-the-art methods. Code is available at <https://github.com/RaccoonDML/CHSNet>.

Index Terms— Crowd counting, noisy annotations

1. INTRODUCTION

Crowd counting is to count the people from a given image in diverse crowded scenes, which is an active computer vision task with a wide range of promising applications in crowd management, traffic monitoring, surveillance systems, etc. Existing methods can be roughly categorized into detection-based [1], count-based [2] and density-map-based [3, 4, 5, 6]. The detection-based methods [1] require laborious annotations (*i.e.* the bounding boxes) to directly detect all persons in an image while the count-based methods [2] only predict the total number of people, suffering from weak supervision.

Unlike the two categories above, density map-based approaches [3, 4] are proposed to estimate the human densities in images, which can balance the performance and annotation cost. Generating a density map only requires point annotations at the center of each head, whose cost is much less



Fig. 1. Noisy annotations in commonly-used datasets. Red points, green points, and yellow arrows denote labeled annotations, missing annotations, and location shifts, respectively.

than the detection-based methods. In addition, density maps can provide more fine-grained pixel-level supervision compared to the count-based methods, which has significantly improved the performance. However, density map-based methods require accurate point annotations to provide reliable pixel-level supervision, which is usually unrealistic because of the potential noises in the labeling process. Fig. 1 shows that missing annotations and location shifts commonly exist among widely-used crowd counting datasets, especially in dense scenes and low-resolution conditions. Therefore, directly using the pixel-level loss function for optimization may compromise the prediction performance. Furthermore, the counting model may memorize the noisy annotations [7].

To address this problem, Ma *et al.* [8] designed a Bayesian loss function for instance-level supervision. Cheng *et al.* [9] proposed a Maximum Excess over Pixels (MEP) loss function, where the region with the maximum loss value is used for optimization. Wan *et al.* [10] modeled the annotation noise using a random variable with Gaussian distribution. Some other methods adopt uncertainty estimation to model the noisy annotations [11, 12]. Recently, Lin *et al.* [13] proposed instance attention loss to exclude the instance in back-propagation if its deviation is too large. However, all of these methods above ignore one critical problem: *how to explore useful supervision in noisy areas?*

†: Corresponding author. jpzhang@fudan.edu.cn. This work was supported in part by the National Natural Science Foundation of China (Nos. 62176059 and 62101136), the Shanghai Municipal Science and Technology Major Project (No. 2018SHZDZX01), the Zhangjiang Laboratory (ZJLab) and the Shanghai Center for Brain Science and Brain-Inspired Technology.

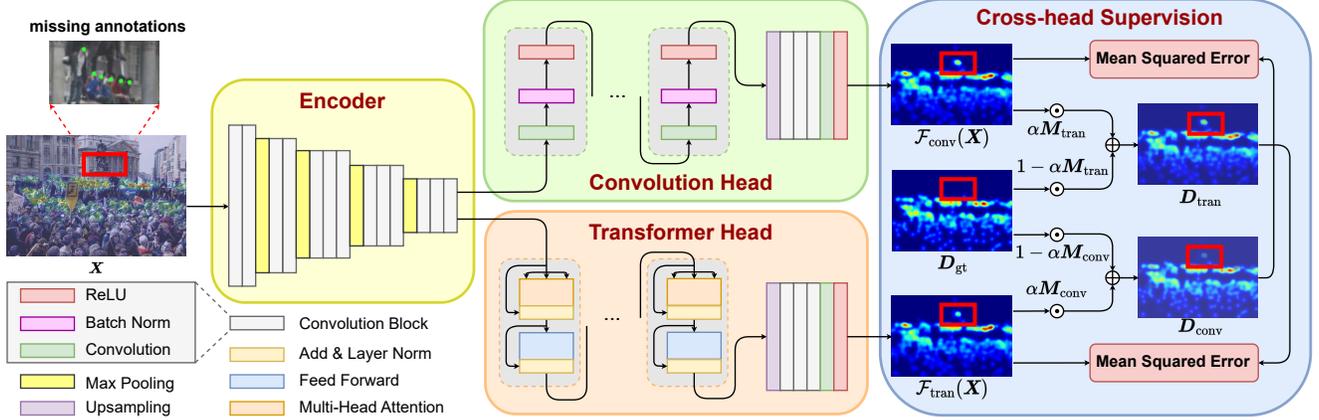


Fig. 2. Architecture of CHS-Net consisting of one shared encoder and two regression heads—one convolution head and one transformer head. At the end of each head, a regression block of the same architecture is attached to produce density map. Red box indicates the area of missing annotations, in which refined supervision is more reliable than ground truth supervision.

Our answer is: *the predicted density map itself can provide useful supervision in noisy areas.* In this paper, we propose to use two regression heads with totally different architecture, *i.e.* one convolution head and one transformer head, to mutually supervise each other in noisy areas. The resultant model, CHS-Net, can synergize different types of inductive biases of convolutions and transformers to boost each other’s performance for better counting. However, the quality of the predicted density map is unsatisfactory in the early stage of training and cannot be directly used for supervision. To make the training reliable, we develop a progressive cross-head supervision learning strategy, that is, the true supervision density maps should be the weighted combinations of the ground truth and predictions from another head, where the weights are linearly increased as the training process goes on.

The main contributions of this work are summarized as follows: 1) we propose a novel model CHS-Net with one convolution head and one transformer head to supervise each other in noisy areas; 2) we design a progressive cross-head supervision learning strategy to make the training process more stable; and 3) our CHS-Net achieves superior performance on several benchmarked datasets.

2. METHODOLOGY

2.1. Network Architecture

Convolutions have strong local modeling ability while transformers [14] can effectively capture global context dependencies. We tackle the noisy annotations present in the data and serve the predictions from two heads as each other’s supervision in noisy areas. Thus, the negative impact of noisy annotations on training is effectively reduced. The inductive biases of the two heads are fully utilized to complement each other.

Fig. 2 presents the proposed CHS-Net consisting of the following components : 1) a shared encoder E that extracts

the features $E(\mathbf{X}) \in \mathbb{R}^{C \times H \times W}$ from the input image \mathbf{X} ; and 2) two regression heads including a convolution head H_{conv} and a transformer head H_{tran} to predict the density maps $H_{\text{conv}}(E(\mathbf{X})) \in \mathbb{R}^{H \times W}$ and $H_{\text{tran}}(E(\mathbf{X})) \in \mathbb{R}^{H \times W}$, where C , H and W are the channels, height, and width of the features maps, respectively. At the end of each head, we adopt an upsampling layer, three convolution blocks, and ReLU activation function to aggregate features and output the density map. The encoder E adopts VGG16 [15] without the last maxpooling layer and fully-connected layers pre-trained on ImageNet [16] for initialization.

On the one hand, the convolution head H_{conv} consists of a series of convolution blocks. Each block is stacked by dilated convolution, batch normalization, and ReLU activation layers, which models the local contextual information. On the other hand, the transformer head H_{tran} contains several transformer layers, whose inputs are the flattened features of $E(\mathbf{X})$ along H and W , to capture global contextual dependencies of current features. For the sake of simplification, we omit the intermediate results and use $\mathcal{F}_{\text{conv}}(\mathbf{X})$ and $\mathcal{F}_{\text{tran}}(\mathbf{X})$ to denote the predictions of two heads.

2.2. Cross-head Supervision

Since the two regression heads generate unreliable density maps at the early stage, we propose a progressive cross-head supervision learning strategy to stabilize the training process. The refined supervision of each head is a weighted combination of another prediction and the ground truth density map. The weights of the complementary head are gradually increased as the training proceeds. Formally, the refined supervision of convolution head \hat{D}_{conv} is defined as follows:

$$\hat{D}_{\text{conv}} = \alpha \mathcal{F}_{\text{tran}}(\mathbf{X}) + (1 - \alpha) D_{\text{gt}}, \quad (1)$$

where α is the combination coefficient to control the importance of two terms and D_{gt} is the ground truth density map

which is usually mislabeled in some area.

Considering that noisy annotations only account for a small part of all annotations, in our method, only a specific mislabeled area uses the refined supervision, while the other areas should use the original ground truth. Typically, the mislabeled examples usually have a large loss [17, 18], as the model would predict the correct labels if it is trained well. Therefore, to select those mislabeled areas, in practice, we first sort the deviation $\epsilon = |\mathcal{F}_{\text{conv}}(\mathbf{X}) - \mathbf{D}_{\text{gt}}| \in \mathbb{R}^{W \times H}$ in descending order and obtain the top $\delta \in [0, 1]$ value as the mask threshold, denoted as t_δ . In a sense, ϵ describes the discrepancy between ground truth and model prediction. Then the selection mask of the convolution head is given by:

$$\mathbf{M}_{\text{conv}} = \mathbb{I}(\epsilon \geq t_\delta) \in \{0, 1\}^{W \times H}, \quad (2)$$

where $\mathbb{I}(\cdot)$ is an indicator function.

Once we have the selection mask, the final supervision for convolution head can be calculated by:

$$\mathbf{D}_{\text{conv}} = \mathbf{M}_{\text{conv}} \odot \widehat{\mathbf{D}}_{\text{conv}} + (1 - \mathbf{M}_{\text{conv}}) \odot \mathbf{D}_{\text{gt}} \quad (3)$$

$$= \alpha \mathbf{M}_{\text{conv}} \odot \mathcal{F}_{\text{tran}}(\mathbf{X}) + (1 - \alpha \mathbf{M}_{\text{conv}}) \odot \mathbf{D}_{\text{gt}}, \quad (4)$$

where \odot represents the element-wise multiplication. Similarly, the final supervision for transformer head \mathbf{D}_{tran} can be calculated using the prediction result by the convolution head. Finally, the overall loss function for optimization is

$$\mathcal{L} = \|\mathcal{F}_{\text{conv}}(\mathbf{X}) - \mathbf{D}_{\text{conv}}\|_2^2 + \|\mathcal{F}_{\text{tran}}(\mathbf{X}) - \mathbf{D}_{\text{tran}}\|_2^2. \quad (5)$$

As a result, the supervision of mislabeled areas is refined from another head.

2.3. Progressive Learning Strategy

The predicted density maps of CHS-Net are unstable in the early stage of training. Therefore, they cannot be directly used for cross-head supervision. To make the early training process stable, we develop a progressive cross-head supervision learning strategy, that is, the noise ratio δ and the combination coefficient α are linearly increased to the preset maximum value as the training process goes on. Formally, the noise ratio and the combination coefficient at the i -epoch are calculated as follows:

$$\delta_i = \delta_{\text{max}} * i/T, \quad \alpha_i = \alpha_{\text{max}} * i/T, \quad (6)$$

where δ_{max} and α_{max} are the predefined maximum noise ratio and maximum combination coefficient, respectively. T denotes the maximum epoch for training.

3. EXPERIMENTS

3.1. Experiment Setups

Datasets. We evaluate our method on three widely-used datasets: ShanghaiTech Part A&B [3] and UCF-QNRF [19].

ShanghaiTech Part A&B contains 482 images (300/182 for training/validation) and 716 images (316/400 for training/validation), respectively. UCF-QNRF includes 1535 high-resolution images (1201/334 for training/validation). This setting covers from sparse scenes to dense scenes and from small dataset to large dataset.

Evaluation metrics. As CHS-Net has two predicted density maps, we use their averaged density map as the final result for evaluation. Mean Absolute Error (MAE) and Mean Squared Error (MSE) are adopted for evaluation.

Table 1. Comparisons with the state-of-the-arts methods on SHA, SHB, and QNRF.

Methods	SHA	SHB	QNRF
	MAE / MSE	MAE / MSE	MAE / MSE
CSRNet [4]	68.2 / 115.0	10.6 / 16.0	- / -
SANet [20]	67.0 / 104.5	8.4 / 13.6	- / -
TEDnet [21]	64.2 / 109.1	8.2 / 12.8	113.0 / 188.0
BL [8]	62.8 / 101.8	7.7 / 12.7	88.7 / 154.8
DM-Count [22]	59.7 / 95.7	7.4 / 11.8	85.6 / 148.3
MCC [23]	71.4 / 110.4	9.6 / 15.0	- / -
NoisyCC [10]	61.9 / 99.6	7.4 / 11.3	85.8 / 150.6
GL [24]	61.3 / 95.4	7.3 / 11.7	84.3 / 147.5
LibraNet [25]	55.9 / 97.1	7.3 / 11.3	88.1 / 143.7
GauNet(CSRNet) [26]	61.2 / 97.8	7.6 / 12.7	84.2 / 152.4
CHS-Net (ours)	<u>59.2 / 97.8</u>	7.1 / 12.1	83.4 / 144.9

Implementation details. We adopt the same data preprocessing in [8]. Ground truth density maps are generated using a fixed Gaussian kernel of size 15. Random scaling, cropping, and horizontal flipping are employed as data augmentation with an image size of 512×512 . Adam optimizer [27] was used with an initial learning rate of 4.0×10^{-5} and weight decay of 1.0×10^{-5} . We use the cosine learning rate scheduler with a maximum epoch of 1,000. For hyperparameters of cross-head supervision, δ_{max} is set to 0.1 for SHA and 0.05 for others. α_{max} is set to 0.5 for QNRF and 1.0 for others.

3.2. Comparison with State-of-the-art Methods

We compare our method with several recent state-of-the-art methods in Table 1. CHS-Net consistently achieves the superior counting performance on all benchmark datasets. For SHA and SHB, our method achieves 59.2 and 7.1 in terms of MAE. For QNRF, CHS-Net improves MAE value of the second best GauNet [26] from 84.2 to 83.4.

The predicted density maps of CHS-Net are visualized in Fig. 3. As can be seen, CHS-Net can predict reliable density maps with high counting accuracy in a wide range of scenes and density levels.

3.3. Ablation Studies

In this section, we perform ablation studies on SHA dataset to evaluate the effectiveness of the proposed CHS-Net.

Ablation study on the models. We first evaluate the model with two heads and cross-head supervision. We have the

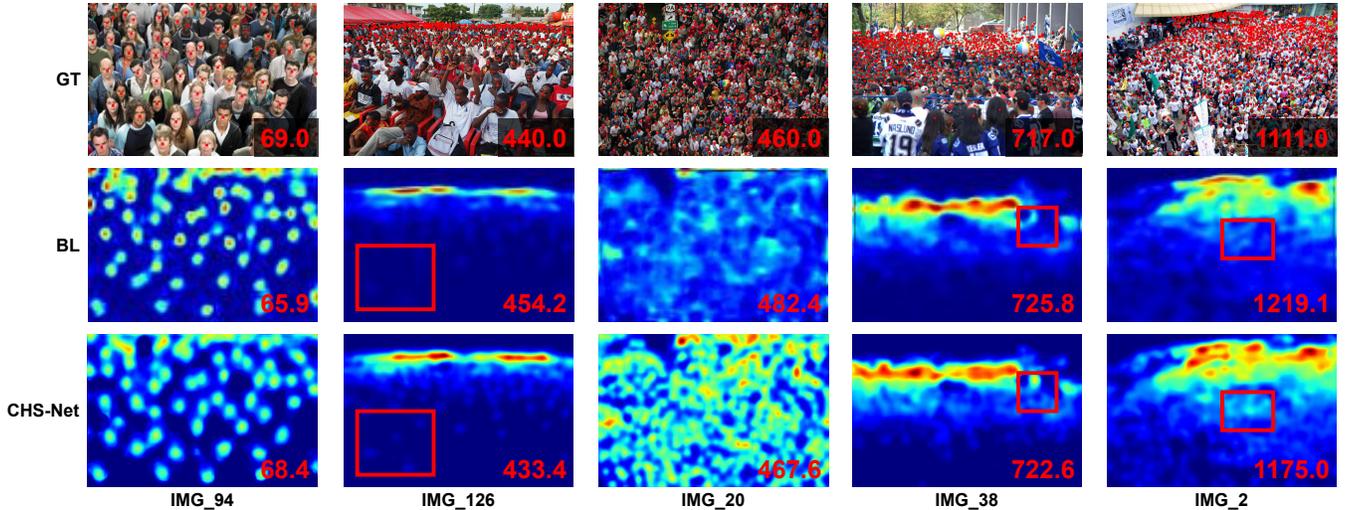


Fig. 3. Visualization on SHA. The first row are the input images. The second and third rows are the predicted density maps by BL [8] and CHS-Net, respectively. Red boxes highlight the differences between them.

Table 2. Ablation study. ‘Average’ represents using the average density map of two heads for evaluation.

Model	Cross-head Supervision	Evaluation Head	MAE	MSE
Conv.	✗	Conv.	63.8	110.4
Tran.	✗	Tran.	62.5	104.1
CHS-Net	✗	Conv.	61.8	107.0
CHS-Net	✗	Tran.	60.7	103.7
CHS-Net	✗	Average	60.7	104.8
CHS-Net	✓	Conv.	59.8	100.4
CHS-Net	✓	Tran.	60.0	96.7
CHS-Net	✓	Average	59.2	97.8

following observations from Table 2: 1) The MAE of transformer head is lower than convolution one by 1.3, which demonstrates the superior global modeling ability of the transformer. 2) Although without cross-head supervision, CHS-Net has achieved significant improvements. 3) With the progressive cross-head supervision learning strategy, the best performance is achieved with MAE of 59.2. 4) We evaluate the performance of every single head in CHS-Net. The average of two heads outperforms any single head in terms of MAE, which further illustrate the advantage of CHS-Net.

Ablation study of two heads on easy/hard samples. Our method is not only an ensemble model with different heads but also contains a self-supervision mechanism implicitly, in which the two heads provide pseudo-labels for each other in the mislabeled area. Based on this idea, we would highlight that the convolution head and transformer head have different learning capabilities. Here, we simply split the hard/easy samples according to the number of humans in an image, *i.e.* the 50% samples with the most humans are hard ones, and we showcase their performance to validate our idea in Table 3. It is found that the convolution head is better at learning easy samples and the transformer head plays the opposite role, respectively. Therefore, the different learning capabilities of

these two heads are of great benefit for providing supervision in the mislabeled area for each other.

Table 3. Ablation study of two heads on easy/hard samples.

MAE / MSE	Easy samples	Hard samples
Conv-head	55.7 / 70.5	216.5 / 299.6
Tran-head	65.2 / 83.4	159.1 / 204.3

Effect of maximum noise ratio. The maximum noise ratio δ_{\max} is an important hyperparameter for CHS-Net. In fact, the maximum noise ratio is a kind of prior knowledge of a specific dataset. We set several values of maximum noise ratio to investigate its effect on model performance. In Table 4, the best performance is obtained when $\delta_{\max} = 0.1$, which means there are nearly 10% noisy annotations in SHA.

Table 4. Effect of maximum noise ratio δ_{\max} .

δ_{\max}	0	0.01	0.05	0.10	0.15	0.30
MAE	60.7	60.8	61.1	59.2	60.5	61.0
MSE	104.8	105.7	99.8	97.8	102.4	105.4

4. CONCLUSION

Noisy annotations are common in crowd counting datasets. To alleviate the negative impact of noisy annotations, we propose CHS-Net, a network with a convolution head and a transformer head to mutually supervise each other in noisy areas. In addition, we develop a progressive cross-head supervision learning strategy to stabilize training process and provide more reliable supervision. Experimental results show the superior performance of our proposed approach. For future work, we will explore more noise robust loss functions to further utilize the ability of CHS-Net. Besides, we may consider to enhance the model sustainable learning ability like [6] so that the noise ratio of training samples from different domains can be adaptively adjusted.

5. REFERENCES

- [1] Vincent Rabaud and Serge Belongie, “Counting crowded moving objects,” in *CVPR*, 2006.
- [2] Antoni B Chan and Nuno Vasconcelos, “Bayesian poisson regression for crowd counting,” in *ICCV*, 2009.
- [3] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma, “Single-image crowd counting via multi-column convolutional neural network,” in *CVPR*, 2016.
- [4] Yuhong Li, Xiaofan Zhang, and Deming Chen, “Csrnet: Dilated convolutional neural networks for understanding the highly congested scenes,” in *CVPR*, 2018.
- [5] Yukun Tian, Yiming Lei, Junping Zhang, and James Z Wang, “Padnet: Pan-density crowd counting,” *IEEE Transactions on Image Processing*, 2019.
- [6] Jiaqi Gao, Jingqi Li, Hongming Shan, Yanyun Qu, James Z Wang, Fei-Yue Wang, and Junping Zhang, “Forget less, count better: a domain-incremental self-distillation learning benchmark for lifelong crowd counting,” *Frontiers of Information Technology & Electronic Engineering*, 2023.
- [7] Sheng Liu, Jonathan Niles-Weed, Narges Razavian, and Carlos Fernandez-Granda, “Early-learning regularization prevents memorization of noisy labels,” in *NeurIPS*, 2020.
- [8] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong, “Bayesian loss for crowd count estimation with point supervision,” in *ICCV*, 2019.
- [9] Zhi-Qi Cheng, Jun-Xiu Li, Qi Dai, Xiao Wu, and Alexander G Hauptmann, “Learning spatial awareness to improve crowd counting,” in *ICCV*, 2019.
- [10] Jia Wan and Antoni Chan, “Modeling noisy annotations for crowd counting,” *NeurIPS*, 2020.
- [11] Min-hwan Oh, Peder Olsen, and Karthikeyan Natesan Ramamurthy, “Crowd counting with decomposed uncertainty,” in *AAAI*, 2020.
- [12] Viresh Ranjan, Boyu Wang, Mubarak Shah, and Minh Hoai, “Uncertainty estimation and sample selection for crowd counting,” in *ACCV*, 2020.
- [13] Hui Lin, Zhiheng Ma, Rongrong Ji, Yaowei Wang, and Xiaopeng Hong, “Boosting crowd counting via multifaceted attention,” in *CVPR*, 2022.
- [14] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, “Attention is all you need,” in *NeurIPS*, 2017.
- [15] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [16] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei, “ImageNet: A large-scale hierarchical image database,” in *CVPR*, 2009.
- [17] Eric Arazo, Diego Ortego, Paul Albert, Noel O’Connor, and Kevin McGuinness, “Unsupervised label noise modeling and loss correction,” in *ICML*, 2019.
- [18] Zhizhong Huang, Junping Zhang, and Hongming Shan, “Twin contrastive learning with noisy labels,” in *CVPR*, 2023.
- [19] 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 *ECCV*, 2018.
- [20] Xinkun Cao, Zhipeng Wang, Yanyun Zhao, and Fei Su, “Scale aggregation network for accurate and efficient crowd counting,” in *ECCV*, 2018.
- [21] 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 *CVPR*, 2019.
- [22] Boyu Wang, Huidong Liu, Dimitris Samaras, and Minh Hoai Nguyen, “Distribution matching for crowd counting,” in *NeurIPS*, 2020.
- [23] Mohsen Zand, Haleh Damirchi, Andrew Farley, Mahdiyar Molahasani, Michael Greenspan, and Ali Etemad, “Multiscale crowd counting and localization by multi-task point supervision,” in *ICASSP*, 2022.
- [24] Jia Wan, Ziquan Liu, and Antoni B Chan, “A generalized loss function for crowd counting and localization,” in *CVPR*, 2021.
- [25] Liang Liu, Hao Lu, Hongwei Zou, Haipeng Xiong, Zhiguo Cao, and Chunhua Shen, “Weighing counts: Sequential crowd counting by reinforcement learning,” in *ECCV*, 2020.
- [26] Zhi-Qi Cheng, Qi Dai, Hong Li, Jingkuan Song, Xiao Wu, and Alexander G Hauptmann, “Rethinking spatial invariance of convolutional networks for object counting,” in *CVPR*, 2022.
- [27] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.