
Boosting Crowd Counting with Transformers

Guolei Sun
CVL, ETH Zurich
Switzerland

Yun Liu
CVL, ETH Zurich
Switzerland

Thomas Probst
CVL, ETH Zurich
Switzerland

Danda Pani Paudel
CVL, ETH Zurich
Switzerland

Nikola Popovic
CVL, ETH Zurich
Switzerland

Luc Van Gool
CVL, ETH Zurich
Switzerland

Abstract

Significant progress on the crowd counting problem has been achieved by integrating larger context into convolutional neural networks (CNNs). This indicates that global scene context is essential, despite the seemingly bottom-up nature of the problem. This may be explained by the fact that context knowledge can adapt and improve local feature extraction to a given scene. In this paper, we therefore investigate the role of global context for crowd counting. Specifically, a pure transformer is used to extract features with global information from overlapping image patches. Inspired by classification, we add a context token to the input sequence, to facilitate information exchange with tokens corresponding to image patches throughout transformer layers. Due to the fact that transformers do not explicitly model the tried-and-true channel-wise interactions, we propose a token-attention module (TAM) to recalibrate encoded features through channel-wise attention informed by the context token. Beyond that, it is adopted to predict the total person count of the image through regression-token module (RTM). Extensive experiments demonstrate that our method achieves state-of-the-art performance on various datasets, including ShanghaiTech, UCF-QNRF, JHU-CROWD++ and NWPU. On the large-scale JHU-CROWD++ dataset, our method improves over the previous best results by 26.9% and 29.9% in terms of MAE and MSE, respectively. Code and models will be available.

1 Introduction

At first sight, counting the size of a crowd present in an image is equivalent to the problem of detecting and counting of person instances [1, 2]. Such direct approaches however have been shown not to perform well, because generic detectors suffer from the small instance size and severe occlusions present in crowded regions [3–5] – typically a person covers only a small number of pixels, and only few body parts are visible (often just the head) [6]. State-of-the-art crowd counting approaches therefore rely on the prediction of crowd density maps, a localized, pixel-wise measure of person presence [7–13, 6, 14, 15, 4, 16–31]. To this end, underlying network architectures need to integrate context across location and scales [4, 12, 32–34]. This is crucial due to the vast variety of possible appearances of a given crowd density. In other words, the ability to integrate a large context makes it possible to adapt the density estimation to an expectation raised by the given scene, beyond the tunnel vision of local estimation. *Geometry* and *semantics* are two of the main aspects of scene context [35–39], that can serve this goal for crowd counting [35, 40]. Unfortunately, even if we manage to model and represent such knowledge, it is very cumbersome to obtain, and therefore not practical for many applications of image-based crowd counting. This also reflects the setup of the most popular crowd counting challenge datasets considered in this paper [8, 6, 41, 42].

On the bright side, even in the absence of such direct knowledge, we can benefit from the recent progress in geometric and semantic learning on a conceptual level – by studying the inductive biases. In fact, the development of computer vision in the last decade demonstrated the possibility to implicitly learn representations capturing rich geometric [43, 44] and semantic [45–47] information from a single image. Recently, the advantageous nature of global interaction over CNNs has been demonstrated for both geometric features for monocular depth prediction [48], as well as for semantic features in segmentation [45]. The aforementioned works attribute the success of the transformer [49, 50] to global receptive fields, which has been a bottleneck in previous CNN-based approaches. Moreover, CNNs by design apply the same operation on all locations, rendering it a sub-optimal choice for exploiting information about the geometric and semantic composition of the scene.

As geometric and semantic understanding are crucial aspects of scene context for the task of crowd counting, we hypothesize that superior capabilities of transformers on these aspects are also indicative of a more suitable inductive bias for crowd counting. To investigate our hypothesis, we adapt state-of-the-art vision transformers [50, 51, 45] for the task of crowd counting.

Unlike image classification [50], crowd counting is a dense prediction task. Following our previous discussion, the learning of crowd counting is also predicated on the global context of the image. To capture both spatial information for dense prediction, as well as the necessary scene context, we maintain both local tokens (representing image patches) and a context token (representing image context). We then introduce a token attention module (TAM) to refine the encoded features informed by the context token. We further guide the learning of the context token by using a regression token module (RTM), that accommodates an auxiliary loss on the regression of the total count of the crowd. Following [45], the refined transformer output is then mapped to the desired crowd density map using two deconvolution layers. Please refer to Fig. 1 for an illustration of the overall framework.

In particular, our proposed TAM module is designed to address the observation that the multi-head self-attention (MHSA) in vision transformers only models spatial interactions, while the tried-and-true channel-wise interactions have also been proved to be of vital effectiveness [52, 53]. To this end, TAM imprints the context token on the local tokens by conditional recalibration of feature channels, therefore explicitly modelling channel-wise interdependencies. Current widely-used methods to achieve this goal includes SENet [52] and CBAM [53]. They use simple aggregation technique such as global average pooling or global maximum pooling on the input features to obtain channel-wise statistics (global abstraction), which are then used to capture channel-wise dependencies. For transformers, we propose a natural and elegant way to model channel relationships by extending the input sequence with a context token and introducing the TAM to recalibrate local tokens through channel-wise attention informed by the context token. The additional attention across feature channels further facilitates the learning of global context.

We also adopt context token which interacts with other patch tokens throughout the transformers to regress the total crowd count of the whole image. This is achieved by the proposed RTM, containing a two-layer MLP. On the one hand, the syzygy of TAM and RTM forces the context token to collect and distribute image-level count estimates from and to all local tokens, leading to a better representation of context token. On the other hand, it helps to learn better underlying features for the task and reduce overfitting within the network, similar to *auxiliary-task learning* [54, 55].

In summary, we provide another perspective on density-supervised crowd counting, through the lens of learning features with global context. Specifically, we introduce a context token tasked with the refinement of local feature tokens through a novel framework of token-attention and regression-token modules. Our framework thereby addresses the shortcomings of CNNs with regards to capturing global context for the problem of crowd counting. Extensive experiments demonstrate that our approach achieves state-of-the-art results on widely-used datasets, including ShanghaiTech, UCF-QNRF, JHU-CROWD++ and NWPU. On JHU-CROWD++, in particular, our model improves over previous best reported result by 26.9% in MAE.

2 Related works

2.1 Object Counting

The goal of object counting is to count the number of target objects in an image. For the case of natural scenes, counting is performed per object class on normal and everyday images [56–60], such

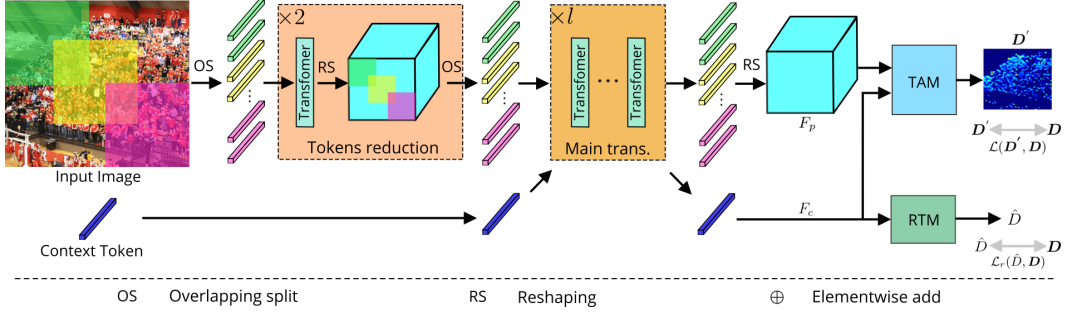


Figure 1: **Network Overview.** The input image is first split into overlapping patches. Then, those patches go through tokens reduction block and main transformer to learn features with global information. To abstract global information, context token (blue vector) is added to the input sequence before the main transformer. The encoded features are processed by token-attention module (TAM) and regression-token module (RTM). The small decoder after TAM is not shown for simplicity.

as images on PASCAL [61] and COCO [62] datasets. For the case of surveillance scenes (*e.g.* crowd counting), counting is conducted in a more constrained setting and a large number of counts for the target class (people) needs to be estimated. Because of its promising application in real applications such as people flow statistics [63], transport safety surveillance [64] and animal conservation [65], it has attracted much attention in computer vision community.

Most crowd counting methods are based on convolutional neural networks (CNNs), which can be divided into three categories: counting by regression [66–69], counting by detection [1, 70, 71, 2], and counting by estimating density maps [7–13, 6, 14, 15, 4, 16–30]. The regression-based methods simply regress the total count of the crowd in the image while the location of the people is not considered. Detection-based approaches first detect the people and then count the number of detections. However, those methods do not perform well since they are not able to estimate the density of the crowd in the image, *i.e.*, the dense region and sparse region. Hence, the mainstream direction of crowd counting is to estimate the density map of the image and then sum over the density map to obtain the total count. All those methods are based on CNNs, which only exploit local context to learn features. To the best of our knowledge, there is only one work [72] adopting vision transformers to conduct crowd counting. However, the method of [72] is concerned with weakly supervised crowd counting in the sense of only regressing the total count, where dot annotations are not available. Its performance therefore cannot compete with the mainstream point-supervised crowd counting methods (using dot labels) on most standard benchmarks [8, 6, 42]. To the best of our knowledge, in this paper we investigate point-supervised crowd counting using vision transformers for the first time and show state-of-the-art performance on popular challenging datasets.

2.2 Vision Transformer

Transformer, relying on self-attention mechanism [49], was first introduced in the domain of natural language processing [49, 73–75] and has been dominating the area ever since. In general, a transformer contains a multi-head self-attention (MHSA) module and a multi-layer perceptron (MLP), to model the contextual information within input sequences through global interaction. Recently, pioneer works such as ViT [50] and DETR [76], proposed to utilize transformers to solve vision problems, by representing images as sequences of patches. It has been shown that transformers are effective in tasks such as image classification [51, 77], object detection [76], semantic/instance segmentation [45], and video segmentation [78]. Specially, ViT [50] proposed to cut the image into patches, which are then converted to sequences of features and used as inputs to the standard transformers. Differently, DETR [76] adopted the features output from CNNs as the input to transformers and the learned features are used for object detection. However, the exploration of transformers for crowd counting has been limited to regression of the total count [72]. In this paper, we demonstrate the power of transformers in point-supervised crowd counting setup, where persons are represented with a binary pixel-wise map.

3 Method

3.1 Problem Definition

Given a training dataset of natural images $\{I_i\} \subseteq \mathbb{R}^{c \times h \times w}$ and corresponding crowd density label maps $\{D_i\} \subseteq \mathbb{R}^{h \times w}$, our goal is to find and learn a neural network model $\mathcal{M} : \mathbb{R}^{c \times h \times w} \rightarrow \mathbb{R}^{h \times w}$, that estimates the crowd density map $D' = \mathcal{M}(I)$ and therefore counts the number of visible people $\|D'\|_1$ from an unseen image I . The sought model \mathcal{M} further minimizes the generalization error of the crowd count estimate on a representative test dataset.

3.2 Transformer-based Crowd Counting

Most conventional crowd counting methods [12, 4, 14] are based on CNNs, and generally consider crowd counting as a dense prediction task. An encoder is used to learn high-level features while a small decoder is adopted to process the encoded output and predict a density map. Since the CNN-based encoder can only exploit the local information within the fix-sized window, some approaches are proposed to increase the receptive fields, by dilated convolution [12, 79] or using deeper network [16]. In this section, we present our transformer-based approach for crowd counting, which is designed to overcome this limitation by explicitly modelling global context. Our presentation follows the data flow of our framework as depicted in Fig. 1.

Overlapping Split. In the popular ViT [50], the input image is split into non-overlapping patches, leading to the problem that the local structure around the patches is destroyed. Instead, we split the input into overlapping patches, following [51]. The process of overlapping split is similar to an convolution operation and the patch size of $k \times k$ is similar to the kernel size. Specifically, the input image I is first padded by p pixels on each side. The overlapping patches are obtained by moving the patch window ($k \times k$) across the whole image with stride s ($s < k$). Each patch has $k \times k \times c$ elements, which are flattened to \mathbb{R}^{ck^2} . The length of patches is given by

$$N_0 = h_0 \times w_0, \tag{1}$$

where $h_0 = \lfloor \frac{h+2p-k}{s} + 1 \rfloor$ and $w_0 = \lfloor \frac{w+2p-k}{s} + 1 \rfloor$. After concatenating all patches together, image tokens are obtained, denoted by $Z_0 \in \mathbb{R}^{N_0 \times ck^2}$. Later, we process Z_0 by tokens reduction block, followed by main transformer.

Tokens Reduction. We first input Z_0 to a transformer layer and obtain Z_1 , formulated as

$$Z_1 = \text{MLP}(\text{MHSA}(Z_0)), \tag{2}$$

where $Z_1 \in \mathbb{R}^{N_0 \times d}$, and d is the dimension of *query*, *key*, and *value*. Since the sequence length N_0 is relatively large due to the overlapping split, we reshape Z_1 back to $\mathbb{R}^{h_0 \times w_0 \times d}$ and perform overlapping split again to reduce the spatial size by stride s . Let Z'_1 be the obtained tokens with size of $\mathbb{R}^{N_1 \times dk^2}$, and $N_1 = h_1 \times w_1$, where $h_1 = \lfloor \frac{h_0+2p-k}{s} + 1 \rfloor$ and $w_1 = \lfloor \frac{w_0+2p-k}{s} + 1 \rfloor$. Following [51], this process is repeated twice and we obtain $Z'_2 \in \mathbb{R}^{N_2 \times dk^2}$, where $N_2 = h_2 \times w_2$. The length of sequence N_2 is thereby reduced to a manageable scale. Since the dependency among those pixels around the original non-overlapping split (as in ViT [50]) is well-modelled, we fix the length of sequence as $N = N_2$ and do not reduce it further, in order to maintain both the representation capability and efficiency. After projecting Z'_2 to $T \in \mathbb{R}^{N \times d}$, we process T by deep-narrow ViT [50], following [51].

Context Token. Recall that we approach crowd counting as a dense prediction problem, and each patch token transforms local RGB input to a local density map prediction. Therefore, even though the patch tokens T are in principle able to interact globally in ViT [50], our mode of dense supervision renders each token to be primarily concerned with its local region. In order to foster global information exchange without compromising capacity for local features, we delegate the collection of global context to a context token t_{con} . This way, patch tokens remain dedicated to the local predictions. The context token is the key input for the TAM module as described in §3.3, which disseminates the global context back to the local tokens. In §3.4 we explain how to guide the learning of context token through the RTM module. But first, we give a brief description of the main transformer of our framework.

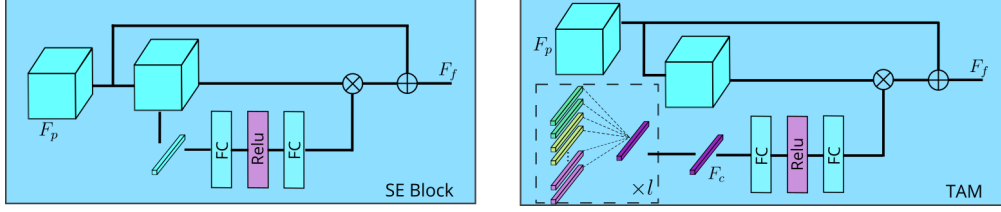


Figure 2: Comparison between SE block and the proposed Token-attention module (TAM).

Main Transformer. The main transformer follow the same architecture as ViT [50], but have less channels in intermediate layers to reduce redundancy within original ViT model. As laid out above, we append the context token t_{con} to the patch tokens T to facilitate global-local interaction. Following [50], position embedding E is also added. The main transformer is denoted as follows

$$\begin{aligned}
 T_0 &= [T; t_{con}] + E, \quad E \in \mathbb{R}^{(N+1) \times d}, \\
 T'_i &= \text{MHSA}(T_{i-1}) + T_{i-1}, \quad i = 1, \dots, l, \\
 T_i &= \text{MLP}(T'_i) + T'_i, \quad i = 1, \dots, l.
 \end{aligned} \tag{3}$$

Here, l is the number of layers in the main transformer. T_l is the feature sequence from the last layer of transformers. It has global receptive fields which are effective for crowd counting task. Since a context token is added in the beginning, we split T_l as follows

$$F_p = T_l[:N], \quad F_c = T_l[N], \tag{4}$$

where $F_p \in \mathbb{R}^{N \times d}$ is the feature corresponding to image patches, and $F_c \in \mathbb{R}^d$ is the feature vector corresponding to context token t_{con} . To recover spatial structure, F_p is reshaped to $\mathbb{R}^{d \times h_2 \times w_2}$. F_p is further refined by token-attention module (TAM) to predict the density map and F_c is used by the proposed regression-token module (RTM) to predict the overall count for the image .

3.3 Token-attention Module (TAM)

The task of the TAM is to refine the local feature map F_p used to predict the crowd density map, conditioned upon the global feature token F_c . To this end, F_p is reshaped from the sequence of tokens (T_l) corresponding to image patches from the last transformer layer. Recall that T_l is produced from T_{l-1} , we detail this process as follows to benefit analysis,

$$\begin{aligned}
 T'_i &= \text{softmax}\left(\frac{T_{l-1} \mathbf{W}_Q (T_{l-1} \mathbf{W}_K)^T}{\sqrt{d}}\right) T_{l-1} \mathbf{W}_V + T_{l-1}, \\
 T_i &= \text{MLP}(T'_i) + T'_i,
 \end{aligned} \tag{5}$$

where $T_{l-1}/T'_i/T_i \in \mathbb{R}^{N \times d}$, and $\mathbf{W}_Q/\mathbf{W}_K/\mathbf{W}_V \in \mathbb{R}^{d \times d}$ are the learnable parameters for generating (*query, key, value*). Here, multi-head self-attention (MHSA) is represented by a special case where a single self-attention (SA) operation is performed for simplicity of notation. As shown in Eq. 5, $T'_i[i]$ corresponding to a specific image patch or context token, is generated by a weighted summation of T_{l-1} . Therefore, transformers are inherently equipped with spatial attention mechanism which is able to pay more attention to the relevant spatial regions (tokens). However, the channel interdependencies are not explicitly modelled in Eq. 5, which could be important for subsequent processing. In fact, this is confirmed by our experiments: while no improvements are obtained by adding spatial attention (CBAM [53]), introducing channel attention through TAM yields better predictions. Explicitly modelling channel relationships, so that the network has the capability to focus on important feature channels, can lead to enhanced features.

Global Abstraction. Global abstraction is used to provide cues for channel interdependencies. In SENet [52], global abstraction is obtained by conducting global average pooling across spatial dimensions on the input itself, while CBAM [53] merges both global average pooling and global maximum pooling. For transformers, we propose a natural and elegant approach to abstract global information, by extending the input sequence with a context token, as introduced above. Since the obtained feature F_c from context token has a global overview on all patches throughout transformer layers, we adopt it to provide information on which channels are important for predicting density map.

The comparison between SE and TAM is shown in Fig. 2, where the sigmoid is omitted for simplicity. The superiority of the proposed TAM over SE [52] and CBAM [53] is validated by experiments (§4.3).

Token-adaptive Recalibration. To capture channel-wise interactions, F_c is projected by a two-layer MLP with ReLU activation, learning a weight vector F'_c which is used to re-weight the feature across the channels. F'_c is obtained by

$$F'_c = \text{sigmoid}(\text{MLP}(F_c)). \quad (6)$$

Here, $F'_c \in \mathbb{R}^d$ and sigmoid function is used to squeeze each element to a range of 0 and 1. We use a convolution layer to prepare F_p for the re-weighting and obtain F'_p . After the recalibration, we add a skip connection [80] with original F_p , and obtain final feature map $F_f \in \mathbb{R}^{h_2 \times w_2 \times d}$, like

$$F_f = F_p + F'_p \otimes F'_c. \quad (7)$$

Through the token-attention module (TAM), the network can increase sensitivity to informative features which are important for downstream processing.

3.4 Regression-token Module (RTM)

Similar to the class token, which is used to find the object class over all patches in ViT [50], context token is used to collect global context over the whole image. It has a global overview of all image patches, through exchanging information with feature vector of each patch throughout all layers. Different from previous work [45], which only uses tokens corresponding to image patches when conducting downstream vision tasks, we adopt F_c to predict the overall count for the whole image. A two-layer MLP with ReLU activation is used to predict the total count \hat{D} , given by

$$\hat{D} = \text{MLP}(F_c). \quad (8)$$

We use L_1 loss to reduce the difference between \hat{D} and ground-truth count, as follows,

$$\mathcal{L}_r(\hat{D}, \mathbf{D}) = |\hat{D} - \|\mathbf{D}\|_1|. \quad (9)$$

Note that we only use this module during training, the predicted count for an image during test is obtained by summing over the predicted density map \mathbf{D}' (§3.5), following other density-map based approaches [12, 30]. The benefits of RTM are two-fold. First, it forces to learn better context-token feature, which provides better information on the importance of each channel and enhance the final feature map F_f . In addition, it helps to learn better underlying feature representations and reduce over-fitting. This can be understood from a view of *auxiliary-task learning* [54, 55], which has shown to be effective in segmentation [81, 47]. By forcing the network to count the crowd using only context-token feature, better information exchange between context token and image patches can be achieved.

3.5 Density Map Prediction and Loss Functions

Following [30], the feature map F_f is processed by a small decoder containing two convolutional layers to predict a density map \mathbf{D}' . We adopt the losses from [30] to supervise density map \mathbf{D}' , which are combined with the loss (Eq. 9) on the context token from regression-token module (RTM). Specifically, the losses for learning density map is a combination of counting loss, optimal transport loss [82] and variation loss, denoted as

$$\mathcal{L}_d(\mathbf{D}', \mathbf{D}) = \|\|\mathbf{D}'\|_1 - \|\mathbf{D}\|_1\| + \lambda_1 \mathcal{L}_{OT}(\mathbf{D}', \mathbf{D}) + \lambda_2 \mathcal{L}_{TV}(\mathbf{D}', \mathbf{D}). \quad (10)$$

The first term of Eq. 10 is to minimize the difference of the total count between the predicted density map and the ground-truth binary mask. The second term (optimal transport loss) and the third term (variation loss) are used to minimize the distribution difference between \mathbf{D}' and \mathbf{D} by regarding the density map as a probability distribution. λ_1 and λ_2 are the weights for the corresponding loss term and we use the same values as in [30] for all our experiments. Please refer to [30] for more details. The total loss function for the proposed method is given by

$$\mathcal{L} = \mathcal{L}_d(\mathbf{D}', \mathbf{D}) + \lambda \mathcal{L}_r(\hat{D}, \mathbf{D}), \quad (11)$$

where λ is weight for $\mathcal{L}_r(\hat{D}, \mathbf{D})$, computed from regression-token module (RTM). The final predicted count for inference is the summation over the predicted density map, given by $\|\mathbf{D}'\|_1$.

4 Experiments

We conduct extensive experiments on four popular benchmark crowd counting datasets [8, 6, 41, 42] to validate the effectiveness of the proposed approach. In this section, we first introduce our experimental setting, followed by state-of-the-art comparisons with previous methods. Finally, we show ablation studies to examine the effectiveness of different components of our model.

4.1 Experimental Setting

Implementation Details. The number of layers (l) in the main transformer is set to be 14. For initialization, we use the official T2T-ViT-14 model [51] pretrained on ImageNet [83]. For data augmentations, we adopt random cropping and random horizontal flipping for all experiments. The initial learning rate and weight decay are set to be $1e-5$ and $1e-4$, respectively. For the optimizer, we use Adam because of its simplicity. Following [30], λ_1 and λ_2 are set to be 0.1 and 0.01, respectively. Following [45], we also compute auxiliary losses at transformer layers of T_5 , T_8 , and T_{11} to provide more supervision during training while only output from last layer is used for prediction. Our method is implemented in the PyTorch framework [84], and experiments are conducted on a single NVIDIA Tesla V100 GPU. To provide full details, we will release our implementations.

Datasets. Experiments are conducted on four challenging datasets: ShanghaiTech [8], UCF-QNRF [6], JHU-CROWD++ [41], and NWPU [42]. Specifically, ShanghaiTech contains 1,198 images with 330,165 annotations, and UCF-QNRF has 1,535 images with more than one million counts. JHU-CROWD++ and NWPU are two largest-scale and most challenging crowd counting benchmarks. JHU-CROWD++ consists of 4,822 images from diverse scenes with more than 1.5 million dot annotations, and NWPU contains 5,109 images with more than two million annotations. The results for the test set are obtained from the evaluation server.

Evaluation Metrics. Following previous works [12, 30, 4], we use mean average error (MAE) and mean square error (MSE) to evaluate the counting performance. For NWPU dataset, we also use mean normalized absolute error (NAE) as evaluation metric. MAE, MSE and NAE are defined as:

$$\text{MAE} = \frac{\sum_{i=1}^N |\|D'_i\|_1 - \|D_i\|_1|}{N}, \quad \text{MSE} = \sqrt{\frac{\sum_{i=1}^N (\|D'_i\|_1 - \|D_i\|_1)^2}{N}}, \quad \text{NAE} = \frac{1}{N} \sum_{i=1}^N \frac{|\|D'_i\|_1 - \|D_i\|_1|}{\|D_i\|_1},$$

where N is the total number of images, and D'_i and D_i are predicted density map and ground-truth binary mask for i^{th} image, respectively.

4.2 Crowd Counting Results

The state-of-the-art comparisons on various datasets are shown in Table 1, Table 2, and Table 3. For all datasets, the proposed method achieves new state-of-the-art performance except ShanghaiTech B, where comparable results are obtained. It shows that our approach is stable across different datasets.

In Table 1, Table 2, and Table 3, we also show the results of the baseline model which adopts the same transformers and loss functions as ours, but without using token-attention module (TAM) and regression-token module (RTM). It shows that simply exploiting transformers provides a strong baseline, compared to previous CNN-based methods. However, our model outperforms the baseline model in all experiments (except MSE on NWPU test), validating the effectiveness of the proposed modules.

In all cases, our method largely outperforms DM-count [30] which uses the same decoder and loss functions as ours to learn the density map. For example, compared with DM-count on ShanghaiTech A [8], our model reduces MAE and MSE from 59.7 to 53.1, and from 95.7 to 82.2, respectively. This demonstrates the strong representation capability of global context features in the task of crowd counting.

On large-scale and challenging benchmarks such as JHU-CROWD++ [41] and NWPU [42], our approach significantly outperforms the previous best results. To be more specific, our method improves BL [25], the best method on JHU-CROWD++ test set, by reducing MAE and MSE from 75.0 to 54.8 and from 299.9 to 208.5, respectively. On NWPU dataset, our approach outperforms DM-count [30], the best method on the NWPU test set, by a margin of 6.4 and 21.7 on MAE and MSE, respectively. Note that the annotations for NWPU test set are not publicly available and the corresponding results are obtained from the evaluation server.

Table 1: Comparison with state-of-the-art methods on ShanghaiTech A [8], ShanghaiTech B [8], and UCF-QNRF [6] datasets. The best and second best results are shown in red and blue, respectively.

Method	Publication	Dot	ShanghaiTech A		ShanghaiTech B		UCF-QNRF	
			MAE	MSE	MAE	MSE	MAE	MSE
Crowd CNN [7]	CVPR15	✓	181.8	277.7	32.0	49.8	-	-
MCNN [8]	CVPR16	✓	110.2	173.2	26.4	41.3	277	426
CMTL [9]	AVSS17	✓	101.3	152.4	20.0	31.1	252	514
Switch CNN [10]	CVPR17	✓	90.4	135.0	21.6	33.4	228	445
IG-CNN [11]	CVPR18	✓	72.5	118.2	13.6	21.1	-	-
CSRNet [12]	CVPR18	✓	68.2	115.0	10.6	16.0	-	-
ic-CNN [13]	ECCV18	✓	68.5	116.2	10.7	16.0	-	-
CL-CNN [6]	ECCV18	✓	-	-	-	-	132	191
SANet [15]	ECCV18	✓	67.0	104.5	8.4	13.6	-	-
CAN [4]	CVPR19	✓	62.3	100.0	7.8	12.2	107	183
SFCN [16]	CVPR19	✓	64.8	107.5	7.6	13.0	102	171
PACNN [17]	CVPR19	✓	62.4	102.0	7.6	11.8	-	-
TEDnet [18]	CVPR19	✓	64.2	109.1	8.2	12.8	113.0	188.0
ANF [19]	ICCV19	✓	63.9	99.4	8.3	13.2	110	174
Wan <i>et al.</i> [20]	ICCV19	✓	64.7	97.1	8.1	13.6	101	176
CFF [21]	ICCV19	✓	65.2	109.4	7.2	12.2	-	-
PGCNet [22]	ICCV19	✓	57.0	86.0	8.8	13.7	-	-
BL [25]	ICCV19	✓	62.8	101.8	7.7	12.7	88.7	154.8
L2R [26]	PAMI19	✓	73.6	112.0	13.7	21.4	124.0	196.0
ASNet [27]	CVPR20	✓	57.7	90.1	-	-	91.5	159.7
LibraNet [28]	ECCV20	✓	55.9	97.1	7.3	11.3	88.1	143.7
Yang <i>et al.</i> [85]	ECCV20	✗	104.6	145.2	12.3	21.2	-	-
NoisyCC [29]	NeurIPS20	✓	61.9	99.6	7.4	11.3	85.8	150.6
DM-Count [30]	NeurIPS20	✓	59.7	95.7	7.4	11.8	85.6	148.3
MATT [86]	PR21	✗	80.1	129.4	11.7	17.5	-	-
Baseline	-	✓	57.3	89.0	7.4	12.2	85.7	150.8
Ours	-	✓	53.1	82.2	7.3	11.5	83.8	143.4

Table 2: Comparison with state-of-the-art methods on the JHU-CROWD++ dataset [41].

Method	Publication	Dot	Val		test	
			MAE	MSE	MAE	MSE
MCNN [8]	CVPR16	✓	160.6	377.7	188.9	483.4
CMTL [9]	AVSS17	✓	138.1	379.5	157.8	490.4
SANet [15]	ECCV18	✓	82.1	272.6	91.1	320.4
CSRNet [12]	CVPR18	✓	72.2	249.9	85.9	309.2
CAN [4]	CVPR19	✓	89.5	239.3	100.1	314.0
SFCN [16]	CVPR19	✓	62.9	247.5	77.5	297.6
BL [25]	ICCV19	✓	59.3	229.2	75.0	299.9
MBTTBF [87]	ICCV19	✓	73.8	256.8	81.8	299.1
CG-DRCN [41]	PAMI20	✓	67.9	262.1	82.3	328.0
Baseline	-	✓	47.6	208.5	58.4	232.7
Ours	-	✓	46.5	198.6	54.8	208.5

4.3 Ablation Studies

Following previous works [12, 4, 27, 28], we conduct ablation studies on ShanghaiTech A [8], to show contributions of key components of our method as well as the effect of hyper-parameter.

TAM and RTM. Table 4a shows that the proposed token-attention module (TAM) and regression-token module (RTM) provide progressive improvements over baseline. Specifically, by adding TAM over baseline, we observe improvement of 2.2 and 2.8 on MAE and MSE, respectively. We also compare the proposed TAM block with SE block [52] and CBAM block [53]. The main difference between TAM and SE/CBAM is that the attention weight for TAM is obtained from context token, while SE/CBAM use feature itself to generate attention weight. As shown in Table 4a, TAM outperforms SE/CBAM, demonstrating that context token contains better information to recalibrate features along channels. The result for CBAM which uses both channel attention and spatial attention shows that additionally adding spatial attention does not help feature learning, since transformers naturally uses spatial attention through multi-head self-attention. By further adding RTM, we obtain the state-of-the-art performance.

Table 3: Comparison with state-of-the-art crowd counting methods on the NWPU dataset [42].

Method	Publication	Dot	Val		test		
			MAE	MSE	MAE	MSE	NAE
MCNN [8]	CVPR16	✓	218.5	700.6	232.5	714.6	1.063
CSRNet [12]	CVPR18	✓	104.8	433.4	121.3	387.8	0.604
CAN [4]	CVPR19	✓	93.5	489.9	106.3	386.5	0.295
SFCN [16]	CVPR19	✓	95.46	608.32	105.7	424.1	0.254
BL [25]	ICCV19	✓	93.64	470.38	105.4	454.2	0.203
KDMG [88]	PAMI20	✓	-	-	100.5	415.5	-
NoisyCC [29]	NeurIPS20	✓	-	-	96.9	534.2	-
DM-Count[30]	NeurIPS20	✓	70.5	357.6	88.4	388.6	0.169
Baseline	-	✓	69.0	314.0	86.6	359.1	0.172
Ours	-	✓	53.0	170.3	82.0	366.9	0.164

Table 4: Ablation study on (a) key components of our method and (b) λ on ShanghaiTech A.

(a)			(b)			
Method	MAE	MSE	Method	λ	MAE	MSE
Baseline	57.3	89.0	Ours	0.01	53.2	83.6
Baseline+TAM	55.1	86.2		0.1	53.1	82.2
Baseline+SE [52]	55.9	87.7		0.2	53.2	82.4
Baseline+CBAM [53]	56.2	90.0		0.5	53.3	82.4
Baseline+TAM+RTM	53.1	82.2				

Effect of λ . Table 4b shows the results when varying λ , which is the weight to control the contribution of \mathcal{L}_r in Eq. 11. It shows that when increasing λ from a small value, the performance improves, because \mathcal{L}_r gradually takes effect and better features are learned. However, when further increasing λ , the performance drops. This is because when λ is a big value, the network focuses less on density map, which leads to reduced performance. For other weight parameters (λ_1 and λ_2), we use the default values as [30] and do not tune them.

5 Conclusion

In this paper, we study the value of global context information in crowd counting with point supervision. We build a strong baseline using transformers to encode features with global receptive fields. Based on that, we proposed two novel module: token-attention module (TAM) and regression-token module (RTM). Extensive experiments are conducted to validate the effectiveness of the proposed method. State-of-the-art performance is achieved on ShanghaiTech, UCF-QNRF, JHU-CROWD++, and NWPU. Unlike most previous methods which are based on CNNs, we provide an alternative perspective for extracting valuable features for point-supervised crowd counting, which could bring new insights to this community.

References

- [1] Weina Ge and Robert T Collins. Marked point processes for crowd counting. In *IEEE CVPR*, pages 2913–2920, 2009. 1, 3
- [2] Tao Zhao and Ramakant Nevatia. Bayesian human segmentation in crowded situations. In *IEEE CVPR*, volume 2, pages 452–459, 2003. 1, 3
- [3] Haoyue Bai and S-H Gary Chan. CNN-based single image crowd counting: Network design, loss function and supervisory signal. *arXiv preprint arXiv:2012.15685*, 2020. 1
- [4] Weizhe Liu, Mathieu Salzmann, and Pascal Fua. Context-aware crowd counting. In *IEEE CVPR*, pages 5099–5108, 2019. 1, 3, 4, 7, 8, 9
- [5] Yutao Hu, Xiaolong Jiang, Xuhui Liu, Baochang Zhang, Jungong Han, Xianbin Cao, and David Doermann. NAS-Count: Counting-by-density with neural architecture search. In *ECCV*, pages 747–766, 2020. 1

- [6] Haroon Idrees, Muhammad 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*, pages 532–546, 2018. 1, 3, 7, 8
- [7] Cong Zhang, Hongsheng Li, Xiaogang Wang, and Xiaokang Yang. Cross-scene crowd counting via deep convolutional neural networks. In *IEEE CVPR*, pages 833–841, 2015. 1, 3, 8
- [8] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. Single-image crowd counting via multi-column convolutional neural network. In *IEEE CVPR*, pages 589–597, 2016. 1, 3, 7, 8, 9
- [9] Vishwanath A. Sindagi and Vishal M. Patel. CNN-based cascaded multi-task learning of high-level prior and density estimation for crowd counting. In *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, pages 1–6, 2017. 8
- [10] Deepak Babu Sam, Shiv Surya, and R Venkatesh Babu. Switching convolutional neural network for crowd counting. In *IEEE CVPR*, pages 4031–4039. IEEE, 2017. 8
- [11] Deepak Babu Sam, Neeraj N Sajjan, R Venkatesh Babu, and Mukundhan Srinivasan. Divide and grow: Capturing huge diversity in crowd images with incrementally growing CNN. In *IEEE CVPR*, pages 3618–3626, 2018. 8
- [12] Yuhong Li, Xiaofan Zhang, and Deming Chen. CSRNet: Dilated convolutional neural networks for understanding the highly congested scenes. In *IEEE CVPR*, pages 1091–1100, 2018. 1, 4, 6, 7, 8, 9
- [13] Viresh Ranjan, Hieu Le, and Minh Hoai. Iterative crowd counting. In *ECCV*, pages 270–285, 2018. 1, 3, 8
- [14] Zenglin Shi, Le Zhang, Yun Liu, Xiaofeng Cao, Yangdong Ye, Ming-Ming Cheng, and Guoyan Zheng. Crowd counting with deep negative correlation learning. In *IEEE CVPR*, pages 5382–5390, 2018. 1, 3, 4
- [15] Xinkun Cao, Zhipeng Wang, Yanyun Zhao, and Fei Su. Scale aggregation network for accurate and efficient crowd counting. In *ECCV*, pages 734–750, 2018. 1, 3, 8
- [16] Qi Wang, Junyu Gao, Wei Lin, and Yuan Yuan. Learning from synthetic data for crowd counting in the wild. In *IEEE CVPR*, pages 8198–8207, 2019. 1, 3, 4, 8, 9
- [17] Miaojing Shi, Zhaohui Yang, Chao Xu, and Qijun Chen. Revisiting perspective information for efficient crowd counting. In *IEEE CVPR*, pages 7279–7288, 2019. 8
- [18] 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 *IEEE CVPR*, pages 6133–6142, 2019. 8
- [19] Anran Zhang, Lei Yue, Jiayi Shen, Fan Zhu, Xiantong Zhen, Xianbin Cao, and Ling Shao. Attentional neural fields for crowd counting. In *IEEE ICCV*, pages 5714–5723, 2019. 8
- [20] Jia Wan and Antoni Chan. Adaptive density map generation for crowd counting. In *IEEE ICCV*, pages 1130–1139, 2019. 8
- [21] Zenglin Shi, Pascal Mettes, and Cees GM Snoek. Counting with focus for free. In *IEEE ICCV*, pages 4200–4209, 2019. 8
- [22] Zhaoyi Yan, Yuchen Yuan, Wangmeng Zuo, Xiao Tan, Yezhen Wang, Shilei Wen, and Errui Ding. Perspective-guided convolution networks for crowd counting. In *IEEE ICCV*, pages 952–961, 2019. 8
- [23] Haipeng Xiong, Hao Lu, Chengxin Liu, Liang Liu, Zhiguo Cao, and Chunhua Shen. From open set to closed set: Counting objects by spatial divide-and-conquer. In *IEEE ICCV*, pages 8362–8371, 2019.
- [24] Lingbo Liu, Zhilin Qiu, Guanbin Li, Shufan Liu, Wanli Ouyang, and Liang Lin. Crowd counting with deep structured scale integration network. In *IEEE ICCV*, pages 1774–1783, 2019.
- [25] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. Bayesian loss for crowd count estimation with point supervision. In *IEEE ICCV*, pages 6142–6151, 2019. 7, 8, 9
- [26] Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. Exploiting unlabeled data in CNNs by self-supervised learning to rank. *IEEE TPAMI*, 41(8):1862–1878, 2019. 8
- [27] Xiaoheng Jiang, Li Zhang, Mingliang Xu, Tianzhu Zhang, Pei Lv, Bing Zhou, Xin Yang, and Yanwei Pang. Attention scaling for crowd counting. In *IEEE CVPR*, pages 4706–4715, 2020. 8

- [28] Liang Liu, Hao Lu, Hongwei Zou, Haipeng Xiong, Zhiguo Cao, and Chunhua Shen. Weighing counts: Sequential crowd counting by reinforcement learning. In *ECCV*, pages 164–181, 2020. 8
- [29] Jia Wan and Antoni Chan. Modeling noisy annotations for crowd counting. In *NeurIPS*, pages 3386–3396, 2020. 8, 9
- [30] Boyu Wang, Huidong Liu, Dimitris Samaras, and Minh Hoai. Distribution matching for crowd counting. In *NeurIPS*, 2020. 3, 6, 7, 8, 9
- [31] Le Zhang, Zenglin Shi, Ming-Ming Cheng, Yun Liu, Jia-Wang Bian, Joey Tianyi Zhou, Guoyan Zheng, and Zeng Zeng. Nonlinear regression via deep negative correlation learning. *IEEE TPAMI*, 2019. 1
- [32] Zhiheng Ma, Xing Wei, Xiaopeng Hong, and Yihong Gong. Learning scales from points: A scale-aware probabilistic model for crowd counting. In *ACM Multimedia*, pages 220–228, 2020. 1
- [33] Qingyu Song, Changan Wang, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Jian Wu, and Jiayi Ma. To choose or to fuse? scale selection for crowd counting. In *AAAI*, 2021.
- [34] Wei Xu, Dingkan Liang, Yixiao Zheng, and Zhanyu Ma. Dilated-scale-aware attention convnet for multi-class object counting. *arXiv preprint arXiv:2012.08149*, 2020. 1
- [35] Dongze Lian, Jing Li, Jia Zheng, Weixin Luo, and Shenghua Gao. Density map regression guided detection network for RGB-D crowd counting and localization. In *IEEE CVPR*, pages 1821–1830, 2019. 1
- [36] Miaogen Ling and Xin Geng. Indoor crowd counting by mixture of Gaussians label distribution learning. *IEEE TIP*, 28(11):5691–5701, 2019.
- [37] James Jerome Gibson and Leonard Carmichael. *The senses considered as perceptual systems*, volume 2. Houghton Mifflin Boston, 1966.
- [38] Krzysztof Kopaczewski, Maciej Szczodrak, Andrzej Czyzewski, and Henryk Krawczyk. A method for counting people attending large public events. *Multimedia Tools and Applications*, 74(12):4289–4301, 2015.
- [39] Timur M. Bagautdinov, Alexandre Alahi, F. Fleuret, P. Fua, and S. Savarese. Social scene understanding: End-to-end multi-person action localization and collective activity recognition. In *IEEE CVPR*, pages 3425–3434, 2017. 1
- [40] Jie He, Xingjiao Wu, Jing Yang, and Wenxin Hu. Cpsnet: Crowd counting via semantic segmentation framework. In *2020 IEEE 32nd International Conference on Tools with Artificial Intelligence (ICTAI)*, pages 1104–1110. IEEE, 2020. 1
- [41] Vishwanath Sindagi, Rajeev Yasarla, and Vishal MM Patel. JHU-CROWD++: Large-scale crowd counting dataset and a benchmark method. *IEEE TPAMI*, 2020. 1, 7, 8
- [42] Qi Wang, Junyu Gao, Wei Lin, and Xuelong Li. NWPU-Crowd: A large-scale benchmark for crowd counting and localization. *IEEE TPAMI*, 2020. 1, 3, 7, 9
- [43] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel J Brostow. Digging into self-supervised monocular depth estimation. In *IEEE ICCV*, pages 3828–3838, 2019. 2
- [44] Jamie Watson, Oisín Mac Aodha, Daniyar Turmukhambetov, Gabriel J Brostow, and Michael Firman. Learning stereo from single images. In *ECCV*, pages 722–740, 2020. 2
- [45] Sixiao Zheng, Jiachen Lu, Hengshuang Zhao, Xiatian Zhu, Zekun Luo, Yabiao Wang, Yanwei Fu, Jianfeng Feng, Tao Xiang, Philip H.S. Torr, and Li Zhang. Rethinking semantic segmentation from a sequence-to-sequence perspective with transformers. In *IEEE CVPR*, 2021. 2, 3, 6, 7
- [46] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, pages 801–818, 2018.
- [47] Guolei Sun, Wenguan Wang, Jifeng Dai, and Luc Van Gool. Mining cross-image semantics for weakly supervised semantic segmentation. In *ECCV*, pages 347–365, 2020. 2, 6
- [48] Guanglei Yang, Hao Tang, Mingli Ding, Nicu Sebe, and Elisa Ricci. Transformers solve the limited receptive field for monocular depth prediction. *arXiv preprint arXiv:2103.12091*, 2021. 2
- [49] 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*, pages 6000–6010, 2017. 2, 3

- [50] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. [2](#), [3](#), [4](#), [5](#), [6](#)
- [51] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token ViT: Training vision transformers from scratch on imagenet. *arXiv preprint arXiv:2101.11986*, 2021. [2](#), [3](#), [4](#), [7](#)
- [52] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *IEEE CVPR*, pages 7132–7141, 2018. [2](#), [5](#), [6](#), [8](#), [9](#)
- [53] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. CBAM: Convolutional block attention module. In *ECCV*, pages 3–19, 2018. [2](#), [5](#), [6](#), [8](#), [9](#)
- [54] Spyros Gidaris, Praveer Singh, and Nikos Komodakis. Unsupervised representation learning by predicting image rotations. In *ICLR*, 2018. [2](#), [6](#)
- [55] Augustus Odena, Christopher Olah, and Jonathon Shlens. Conditional image synthesis with auxiliary classifier GANs. In *ICML*, pages 2642–2651. PMLR, 2017. [2](#), [6](#)
- [56] Prithvijit Chattopadhyay, Ramakrishna Vedantam, Ramprasaath R Selvaraju, Dhruv Batra, and Devi Parikh. Counting everyday objects in everyday scenes. In *IEEE CVPR*, pages 1135–1144, 2017. [2](#)
- [57] Hisham Cholakkal, Guolei Sun, Fahad Shahbaz Khan, and Ling Shao. Object counting and instance segmentation with image-level supervision. In *IEEE CVPR*, pages 12397–12405, 2019.
- [58] Issam H Laradji, Negar Rostamzadeh, Pedro O Pinheiro, David Vazquez, and Mark Schmidt. Where are the blobs: Counting by localization with point supervision. In *ECCV*, pages 547–562, 2018.
- [59] Meng-Ru Hsieh, Yen-Liang Lin, and Winston H Hsu. Drone-based object counting by spatially regularized regional proposal network. In *IEEE ICCV*, pages 4145–4153, 2017.
- [60] Hisham Cholakkal, Guolei Sun, Salman Khan, Fahad Shahbaz Khan, Ling Shao, and Luc Van Gool. Towards partial supervision for generic object counting in natural scenes. *IEEE TPAMI*, 2020. [2](#)
- [61] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The PASCAL visual object classes challenge: A retrospective. *IJCV*, 111(1):98–136, 2015. [3](#)
- [62] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. [3](#)
- [63] Jin Zhang, Sheng Chen, Sen Tian, Wenan Gong, Guoshan Cai, and Ying Wang. A crowd counting framework combining with crowd location. *Journal of Advanced Transportation*, 2021, 2021. [3](#)
- [64] Vishwanath A Sindagi and Vishal M Patel. A survey of recent advances in CNN-based single image crowd counting and density estimation. *Pattern Recognition Letters*, 107:3–16, 2018. [3](#)
- [65] Mohammad Sadegh Norouzzadeh, Anh Nguyen, Margaret Kosmala, Alexandra Swanson, Meredith S Palmer, Craig Packer, and Jeff Clune. Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proceedings of the National Academy of Sciences*, 115(25): E5716–E5725, 2018. [3](#)
- [66] Antoni B Chan and Nuno Vasconcelos. Bayesian poisson regression for crowd counting. In *IEEE ICCV*, pages 545–551. IEEE, 2009. [3](#)
- [67] Ke Chen, Chen Change Loy, Shaogang Gong, and Tony Xiang. Feature mining for localised crowd counting. In *BMVC*, page 3, 2012.
- [68] Ke Chen, Shaogang Gong, Tao Xiang, and Chen Change Loy. Cumulative attribute space for age and crowd density estimation. In *IEEE CVPR*, pages 2467–2474, 2013.
- [69] Chuan Wang, Hua Zhang, Liang Yang, Si Liu, and Xiaochun Cao. Deep people counting in extremely dense crowds. In *ACM Multimedia*, pages 1299–1302, 2015. [3](#)
- [70] Min Li, Zhaoxiang Zhang, Kaiqi Huang, and Tieniu Tan. Estimating the number of people in crowded scenes by mid based foreground segmentation and head-shoulder detection. In *IEEE ICPR*, pages 1–4, 2008. [3](#)

- [71] Sheng-Fuu Lin, Jaw-Yeh Chen, and Hung-Xin Chao. Estimation of number of people in crowded scenes using perspective transformation. *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, 31(6):645–654, 2001. 3
- [72] Dingkan Liang, Xiwu Chen, Wei Xu, Yu Zhou, and Xiang Bai. TransCrowd: Weakly-supervised crowd counting with transformer. *arXiv preprint arXiv:2104.09116*, 2021. 3
- [73] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime Carbonell, Quoc V Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2978–2988, 2019. 3
- [74] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, pages 4171–4186, 2019.
- [75] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, and Quoc V Le. XLNet: Generalized autoregressive pretraining for language understanding. In *NeurIPS*, pages 5753–5763, 2019. 3
- [76] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229. Springer, 2020. 3
- [77] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. *arXiv preprint arXiv:2102.12122*, 2021. 3
- [78] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *IEEE CVPR*, 2021. 3
- [79] Shuai Bai, Zhiqun He, Yu Qiao, Hanzhe Hu, Wei Wu, and Junjie Yan. Adaptive dilated network with self-correction supervision for counting. In *IEEE CVPR*, pages 4594–4603, 2020. 4
- [80] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE CVPR*, pages 770–778, 2016. 6
- [81] Zhi Tian, Chunhua Shen, and Hao Chen. Conditional convolutions for instance segmentation. In *ECCV*, pages 282–298, 2020. 6
- [82] Cédric Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008. 6
- [83] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015. 7
- [84] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8026–8037, 2019. 7
- [85] Yifan Yang, Guorong Li, Zhe Wu, Li Su, Qingming Huang, and Nicu Sebe. Weakly-supervised crowd counting learns from sorting rather than locations. In *ECCV*, 2020. 8
- [86] Yinjie Lei, Yan Liu, Pingping Zhang, and Lingqiao Liu. Towards using count-level weak supervision for crowd counting. *Pattern Recognition*, 109:107616, 2021. 8
- [87] Vishwanath A Sindagi and Vishal M Patel. Multi-level bottom-top and top-bottom feature fusion for crowd counting. In *IEEE ICCV*, pages 1002–1012, 2019. 8
- [88] Jia Wan, Qingzhong Wang, and Antoni B Chan. Kernel-based density map generation for dense object counting. *IEEE TPAMI*, 2020. 9