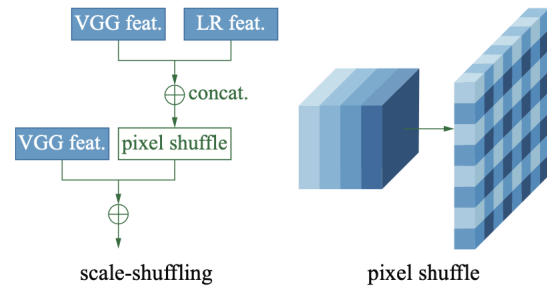
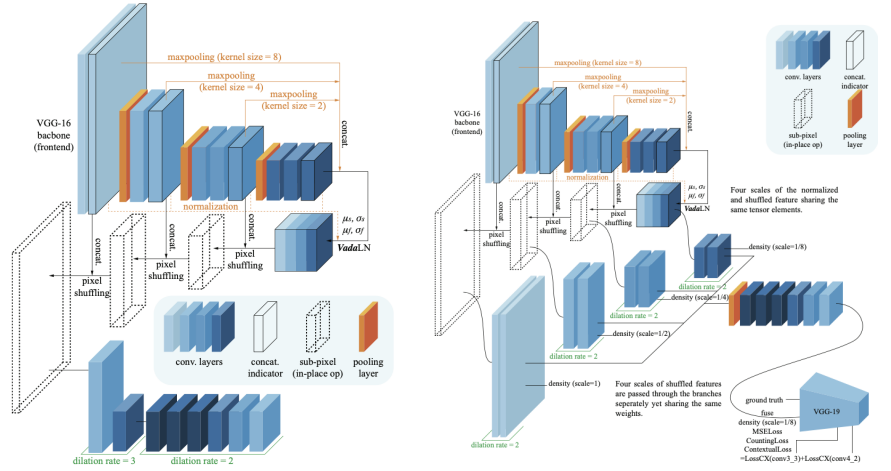


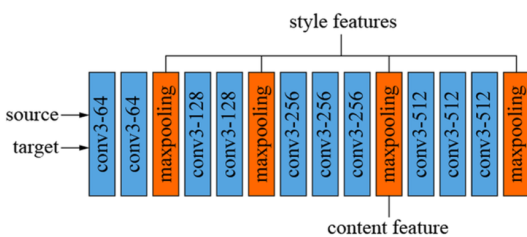
Crowd counting is widely used in real-time congestion monitoring and public security. Due to the limited data, many methods have little ability to be generalized because the differences between feature domains are not taken into consideration. We propose VGG-embedded adaptive layer normalization (VadaLN) to filter the features that irrelevant to the counting tasks in order that the counting results should not be affected by the image quality, color or illumination. VadaLN is implemented on the pretrained VGG-16 backbone. There is no additional learning parameters required through our method. VadaLN incorporates the proposed scale-shuffling modules (SSM) to relax the distortions in upsampling operations. Besides, non-aligned training methodology for the estimation of density maps is leveraged by an adversarial contextual loss (ACL) to improve the counting performance.

We propose VGG-embedded adaptive layer normalization (VadaLN) to address this issue. Moreover, VadaLN also plays a role as a negative-feature filter to generalize the model to handle the images captured in the wild. Adaptive affine transformations using the specific layers in VGG-16 backbone is proposed to adapt LN to arbitrarily given cases of crowd counting. We formulate VadaLN to filter the image styles while keep the content feature for density estimation. Intuitively, different cases contain different illumination styles. Considering of the style-representation layers in pretrained VGG-16, in order to eliminate the influence from such layers while preserving the content of the crowd, the means and variances of these layers are precomputed for the normalization of the content layers.



Shuffle to the Desired Scale: With the operation mentioned above, to the specific task of crowd counting, we firstly resize the input images. Thus, the down-scaled feature sizes in the ReLU layers of the VGG backbone can match those of the sub-pixel convoluted features.

Contextual loss is used to tackle the optimization problems for non-aligned data. The scheme is utilized for those source and target that are not necessarily aligned. With the ground truth generation by the Gaussian kernel, the pixel-wise difference between annotated maps can be considerably high. Thus, simply using pixel-wise loss may lead to the overfitting problem.



$$CX(x, y) = CX(X, Y) = \frac{1}{N} \sum_j \max_i CX_{ij}$$

$$L_{CX}(x, y, l) = -\log(CX(\Phi^l(x), \Phi^l(y)))$$

$$ACL(x, y) = -\lambda_1 \log(CX(\Phi^{\text{conv}_{3.3}}(x), \Phi^{\text{conv}_{3.3}}(y))) \\ -\lambda_2 \log(CX(\Phi^{\text{conv}_{4.2}}(x), \Phi^{\text{conv}_{4.2}}(y))) \\ +\lambda_3 \mathbb{E}_{x \sim p_{\text{pred}}} [\log D_{ACL}^*(x)] \\ +\lambda_3 \mathbb{E}_{x \sim p_{gt}} [\log(1 - D_{ACL}^*(x))]$$