

---

# Cross-Region Building Counting in Satellite Imagery using Counting Consistency

Muaaz Zakria<sup>1</sup> · Hamza Rawal<sup>1</sup> · Waqas Sultani<sup>1</sup> · Mohsen Ali<sup>1</sup>

**Abstract** Estimating the number of buildings in any geographical region is a vital component of urban analysis, disaster management, and public policy decision. Deep learning methods for building localization and counting in satellite imagery, can serve as a viable and cheap alternative. However, these algorithms suffer performance degradation when applied to the regions on which they have not been trained. Current large datasets mostly cover the developed regions and collecting such datasets for every region is a costly, time-consuming, and difficult endeavor. In this paper, we propose an unsupervised domain adaptation method for counting buildings where we use a labeled source domain (developed regions) and adapt the trained model on an unlabeled target domain (developing regions). We initially align distribution maps across domains by aligning the output space distribution through adversarial loss. We then exploit counting consistency constraints, within-image count consistency, and across-image count consistency, to decrease the domain shift. Within-image consistency enforces that the building count in the whole image should be greater than or equal to the count in any of its sub-image. Across-image consistency constraint enforces that if an image contains considerably more buildings than the other image, then their sub-images shall also have the same order. These two constraints encourage the behavior to be consistent across and within the images, regardless of the scale. To evaluate the performance of our proposed approach, we collected and annotated a large-scale dataset consisting of challenging South Asian regions having higher building densities and irregular structures as compared to existing datasets. We perform extensive experiments to verify the efficacy of our approach and report improvements of approximately 7% to

20% over the competitive baseline methods. The dataset and code are available here: <https://github.com/intelligentMachines-ITU/domain-Adaptive-Building-Counting>.

## 1 Introduction

The precise and accurate estimation of the number of buildings is vital for many tasks, such as monitoring economic well-being [53], planning aid for a natural or a man-made disaster-stricken region [1], analyzing poverty [12], and predicting the vitality of a city [37]. Over the years, there has been an effort to approximate the population size by estimating the number of buildings in an area through satellite imagery and through land use/cover data [48]. Building counting acts as an indicator of important metrics such as population density [32,6], and power usage [16]. This information is generally collected through various censuses and surveys or their fusion, requiring costly, expansive, and time-consuming efforts. One way to tackle this challenge is to make the process cost-effective and labor-saving by using deep learning-based methods on satellite images to get an automatic estimate of building counts [38] or through the extraction of buildings [51].

The problem of counting objects in an image has been studied extensively in the last few years. The ubiquitous nature of the counting problem is exhibited in the variety of research works that deal with counting cells in the petri dish [29], estimating crowd size [25], and counting the buildings [38]. In the literature, counting has been performed mainly through segmentation [18], clustering [33], and regression-based methods [25,57,36,45,31]. As compared to counting in normal images, counting the number of buildings (or in short, building counting) from satellite imagery is less explored [35].

The training data on which a deep learning model is trained is called the source domain (data) and while the one

---

E-mail: waqas.sultani@itu.edu.pk ·

<sup>1</sup>Intelligent Machines Lab, Information Technology University, Lahore



**Fig. 1** The left block shows the images captured from colder regions and have sloped roofs while the images in the right block show flat roofs as those areas belong to hotter regions.

on which it is to be tested is called the target domain (data). Ideally, source and target domain data distributions should be the same, however, in real-world applications, source, and target domain distributions are different. This is more apparent in the case of the building counting problem, since depending upon the region, climate, and culture, the building structures are different from each other. In Figure 1, for example, we show images from the colder region having buildings with sloped roofs, whereas the other ones have flat roofs. Similarly density of the neighborhood, building material, time of image acquisition, quality, and resolution of the aerial/satellite imagery, can affect the robustness of the building count. The gap between the distributions of two (source and target) domains, called domain shift, is the reason for this failure. It is well known that the deep learning-based methods, even when trained on large datasets fail to generalize to the new domain [9, 2, 7, 10, 41].

To overcome the limitation due to domain shift, domain adaptation strategies have been applied to numerous problems including crowd counting [7, 2, 41, 9, 10]. Domain adaptation in building counting is a relatively unexplored problem and more focus of these domain adaptation problems is on semantic segmentation or object detection etc [10].

Recently various remote sensing datasets have been introduced for deep learning applications [50, 46]. However, these datasets majorly cover regions from the developed world. When a deep learning model trained on developed regions, is tested on developing or under-developed regions, it declines in performance due to domain shift.

In this paper, we propose to tackle this performance decline in building counting across regions (from developing to under-developed regions). In our case, the source dataset contains regions from developed countries and the target dataset contains developing regions. *Since we do not assume the availability of ground truth data of the target training dataset, we call our approach an unsupervised domain adaptation method.* Given the satellite image, we first automatically obtain building density maps employing [25]. The total building count in an image is produced by summing the whole density map. Learning to predict the map forces the model to learn to localize the buildings it is trying to count. However, due to the domain shift, the density maps predicted on the target domain lack structure. To address this, we propose to use adversarial learning which

forces the model to learn to produce density maps that are indistinguishable across the source and target domains. Furthermore, we design a problem-specific strategy to align two domains for improved building counting. Specifically, we propose two counting consistency constraints, within-image count consistency, and across-image count consistency, to help decrease the domain shift in an unsupervised way. Both of these consistencies should naturally occur in any area. Within-image consistency encompasses the logic that the number of buildings in the whole image should be greater than or equal to the number of buildings in any of the sub-region captured in the image. Across-image consistency constraint on the other hand enforces that if a region contains considerably more number of buildings than other regions, then a considerably large sub-region of the former one will also have more number of buildings than the sub-region of the latter one. These two constraints encourage the behavior to be consistent across and within the images, regardless of the scale.

To summarize, the following are the contributions of our work:

- We attempt to address a new problem of cross-region building counting and localization.
- We propose two problem-specific constraints, count consistency, to direct the unsupervised domain adaptation process. These constraints i.e., the within-image and across-image count consistency constraints force the model to learn a generalized representation of buildings.
- We collect a new large-scale and challenging dataset for building localization and counting with a focus on South Asian regions having irregular building structures.
- We perform extensive experiments to show the effectiveness of our proposed approach and the efficacy of the collected dataset.

We will cover the related works in section 2, explore the used datasets in section 3, explain our adopted methodologies in section 4, detail out our implementation and discuss our results in section 5, and finally conclude in section 7.

## 2 Related Work

In our proposed approach, we tackle the problem of cross-region building counting and localization using counting con-

sistency constraints based on ranking loss and introduce a new dataset. Hence, below we discuss the works related to object counting, ranking, domain adaptation, and remote sensing datasets.

**Object Counting:** Counting objects of interest through computer vision has been done in several application areas which include counting people, animals, fruits, buildings, etc. Counting crowd from images was performed in the early days using detection based methods [49,52,18]. Clustering-based methods have also been employed to count people in crowded scenes [33,45]. However, regression-based counting methods [25,31,36] have generally produced more state-of-the-art results than the other mentioned methods. Recently, Liu et al. [25] have predicted crowd density by encoding the contextual information contained within various scales. To account for the fact that the appropriate scale varies over the image, Kang et al. [13] proposed to weight the generated density maps differently at different scales. Whereas, Li et al. [17] performed crowd counting using multi-resolution context and image quality assessment-guided training. To count building from satellite imagery, Shakeel et al., [38] proposed a regression model using attention-based re-weighting.

Detecting and counting fruits is of unmatched importance in agriculture [34,27,54]. Rahnemoonfar et al. [34] and Liu et al. [27] used deep learning-based methodologies to count fruits in images. Similarly, Zabawa et al. [54] used convolutional neural networks to perform semantic segmentation to detect grapevine berries in images and then count them using a connected component algorithm. The problem of counting vehicles has been tackled using various methodologies in literature [5,31,56].

**Ranking:** The ranking makes sure that the ordering of a list of items is in the correct order. The framework of ranking has been applied to solve various problems of computer vision including improving the counting of crowds in congested scenes using a ranking loss [26], anomaly detection in surveillance videos where the ranking loss was used to localize the anomalies during the training [42] and detection of features by employing a deep ranking framework [47].

**Domain Adaptation:** Domain adaptation is the process of minimizing the effects of a domain shift that arises when training and testing data is drawn from different distributions. Domain adaptation has been used to solve several problems such as object classification [55], detection [58,39,40], semantic segmentation [43,9,2,10], person re-identification [8] and crowd counting [7,21]. Domain adaptation has been addressed by [30] by proposing a new latent sub-domain discovery model for dividing the target domain into sub-domains by considering them a cluster while bridging the domain gap. A weakly supervised domain adaptation network for latent space and output space has been proposed by [10] to diminish the cross-domain gap in satellite and aerial imagery for performing semantic segmentation of built-up

areas. Hossain et al., [7] used domain adaptation for crowd counting where they used semi-supervised domain adaptation using a limited number of labeled images from target data. Finally, they have minimized maximum mean discrepancy (MMD) loss between the generated density maps of source and target. In addition to semi or weakly supervised domain adaptation, some recent works also address unsupervised domain adaptation [22,23,8,59]. Moreover, domain adaptation in remote sensing image classification was also addressed by Zhang et al. [55] and Liu et al. [24] using unsupervised transfer learning. In addition to above cited works, domain shift problems have also been addressed in recommendation systems by [19] using context-aware bandits, by [20] using Collaborative Filtering Bandits, and by similar bandits [14,28].

**Datasets for Remote Sensing:** Several datasets have been introduced for remote sensing applications such as object detection or built-areas detection. One of the most popular datasets for object detection is the xView dataset [15]. It contains bounding box annotations of objects of multifarious classes. The xView dataset consists of a total of 1 million object instances that come under 60 classes. It covers a land area of  $1415km^2$ . A building detection dataset was released by the SpaceNet [46], covering the areas of Rio De Janeiro, Las Vegas, Paris, Shanghai, and Khartoum. A semantic labeling benchmark dataset was launched by ISPRS [11] which contained 2D semantic labels in high quality of two cities in Germany. For the problem of counting built structures from satellite imagery, a large and diverse dataset was introduced by [38] covering urban, hilly, and desert regions. Another important dataset was proposed by [1] to detect destructed sites due to natural and man-made disasters from satellite imagery. Similarly, for object detection in aerial imagery, a large-scale dataset, DOTA, was put forward by [50].

In contrast to the above-mentioned works, our approach focuses on counting buildings in satellite imaging across different regions. Instead of employing within-image counts ranking for self-training [26], we have used across and within-image counts ranking for unsupervised domain adaptation across different regions. Furthermore, we have introduced a new dataset and annotations to demonstrate the efficacy of the proposed approach.

### 3 Dataset Preparation

**xView:** We appropriated xView [15] for the building counting task utilizing it as a *source* dataset. The motivation behind using this dataset is that it is distributed across various parts of the world to detect objects ranging across 60

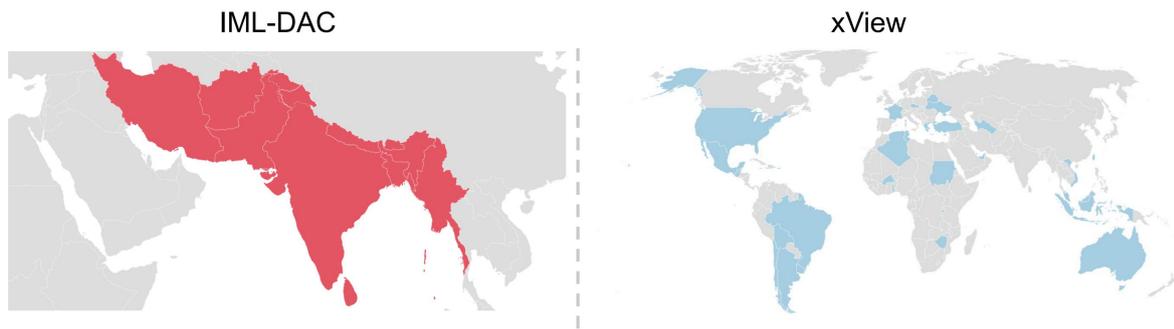


Fig. 2 Side-by-side comparison of geographical locations of IML-DAC and xView datasets.

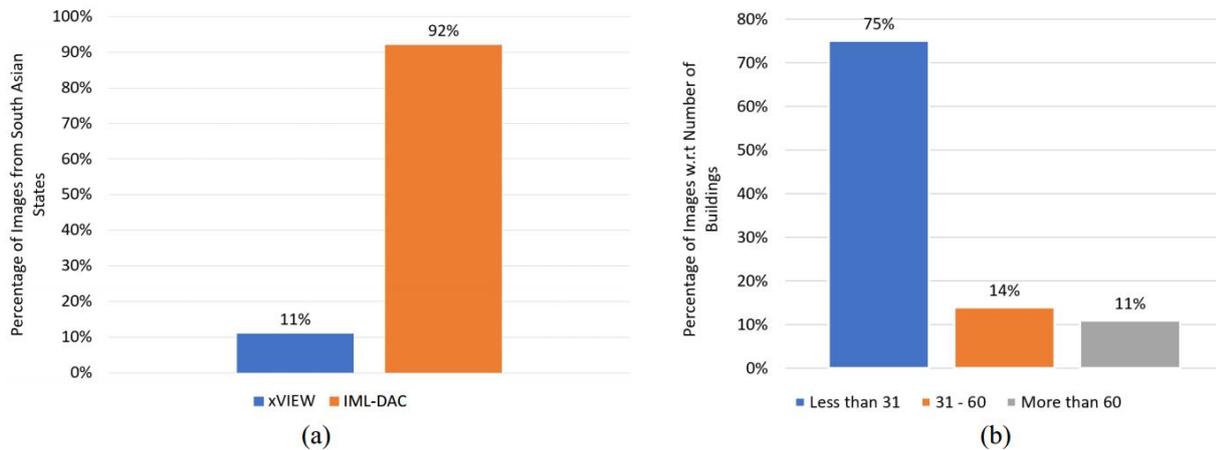


Fig. 3 (a) compares the percentage of South Asian regions contained in xView and IML-DAC. Our dataset contains 81% more images from South Asian regions than xView dataset. (b) shows the distribution of images of xView and IML-DAC datasets with respect to the number of buildings contained in them.

classes which also includes buildings.<sup>1</sup> The xView dataset originally contained 847 training images along with ground truth. The ground truth is available in geoJSON format. The fields in it contained the bounding box label ID, a unique ID for image strips, image filename, coordinates of bounding boxes, and longitude-latitude information of bounding boxes. To generate our desired annotations, we first selected the bounding boxes which covered **Buildings** only. Afterward, we use original ground truth coordinates to generate the points which are in the center of these buildings. After that, we generate non-overlapping patches of size  $500 \times 500$  pixels and selected a total of 4935 patches. The division of these patches with respect to the number of buildings contained in them is shown in Figure 3(b). Note the xView is geographically biased towards developed regions/countries. There are a few images from developing countries that we did not use during training the model.

**IML-DAC:** To evaluate the accuracy of the proposed approach in cross-region building counting, in our experiments,

<sup>1</sup> We have not used DOTA [50] since it does not contain buildings class.

we use the xView as a source (train) dataset and IML-DAC and South Asian regions of xView as target (test) datasets. The geographical locations of areas from which xView and IML-DAC are collected are presented in Figure 2. Figure 3(a) compares the percentage of South Asian regions contained in both datasets and Figure 3(b) shows the distribution of images of xView and IML-DAC with respect to the number of buildings contained in them. Figure 5 communicates a better understanding of the distribution of our dataset with respect to the percentage of images from each region.

In Figure 4, we demonstrate a comparison of both datasets with respect to their counts and structures. Figure 4(a) shows the images (side by side) of xView and IML-DAC having similar building counts. It can be observed that buildings in xView are well-placed and distant while there is no proper planning for building placements in the IML-DAC dataset. Figure 4(b) highlights the difference in structures of both datasets. Most images in xView contain tall buildings while in IML-DAC, the majority of buildings are small in size or either built of non-concrete material. Since our main goal is to count buildings, the images in IML-DAC are annotated through a dotted annotation on each building. Some typical



**Fig. 4** Comparison of images from both datasets. (a) shows the images of xView and IML-DAC for similar building counts. Similarly (b) highlights the difference in building structures. In the last row, we show the annotated points on the buildings. The images are histogram equalized.

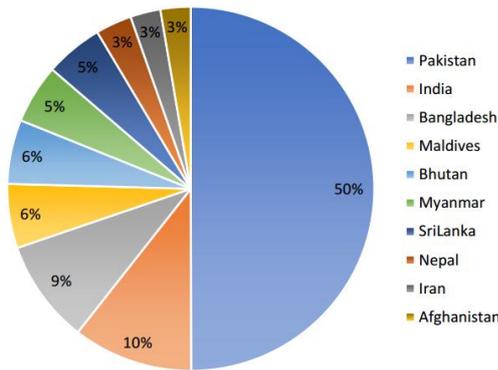
examples of annotations are shown in the bottom row of Figure 4. Note that as mentioned in the ‘xView’ section, dotted annotations for xView are extracted from the bounding box annotations provided by the original authors of xView.

## 4 Methodology

In the following section, we provide the details of each component used in our pipelines and explain all the design choices behind them.

### 4.1 Preliminaries

Let us define the source and target datasets to be  $\mathcal{D}^s = \{(I_i^s, l_i^s), i = 1 \dots N^s\}$  and  $\mathcal{D}^t = \{(I_j^t), j = 1 \dots N^t\}$ . Where  $I_i^s$  and  $I_j^t$  are the satellite imagery patches from the source and target datasets, and  $N^s$  and  $N^t$  are the total number of source and target data image patches respectively. Note that in this paper, we use satellite images from developed regions as a *source* dataset and satellite images from developing regions as *target* dataset. For each image patch



**Fig. 5** Distribution of IML-DAC dataset. The majority of images have been collected from Pakistan across its various cities.

$I_i^s$  in the source domain dataset, we have a ground-truth list  $l_i^s$  of locations where the buildings are present. Using [25] the ground-truth density map  $D_i^s$  is created for each source sample, such that a Gaussian is centered on each location in  $l_i^s$  and the variance of the Gaussian depends upon how far away the other buildings are from the current one.

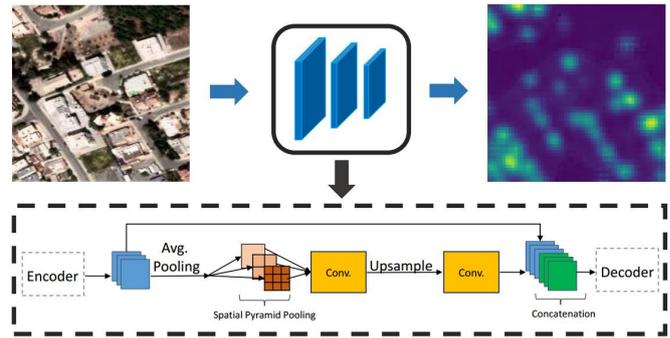
Since the proposed approach works on the building density maps, we use the Context-Aware Convolutional Network (CACN) [25] based on counting pipeline  $f_c$ . Originally this network was introduced to count people in images of crowded scenes. The motivation behind using CACN for counting buildings is the fact that it can encode the contextual information contained within multiple scales in an adaptive manner by incorporating spatial pyramid pooling. Such spatial and contextual variations are also visible in the building dataset. CACN is trained over the source domain using the ground-truth density maps  $D_i^s$ 's.

$$\mathcal{L}_{MSE}^s = \frac{1}{2N^s} \sum_{i=1}^{N^s} \|D_i^s - f_c(I_i^s)\|_2^2, \quad (1)$$

where  $\mathcal{L}_{MSE}^s$  represents the mean square loss on the source dataset,  $N^s$  is the batch size chosen for training,  $D_i^s$  is the ground truth density map and  $f_c(I_i^s)$  is the predicted density map of image patch  $I_i^s$ . Figure 6 shows the example of a building density map generated by CACN. During inference, the predicted count is the summation of the density map. This source-only model serves as a baseline upon which we add our proposed modules to conduct progressive performance enhancement.

## 4.2 Unsupervised Domain Alignment

The source-only trained model from the previous section performs poorly on the target domain. This decrease in performance is attributed to the domain gap between the source and target domains. We design an unsupervised domain adaptation strategy guided by adversarial feature alignment and



**Fig. 6** Generating density maps of buildings from satellite images using Context-Aware Convolutional Network (CACN) [25].

the proposed consistency constraints. The resultant model is robust to domain shift and is more generalized than the baseline source-only model.

### 4.2.1 Distribution Map Alignment (DMA) using Adversarial Learning

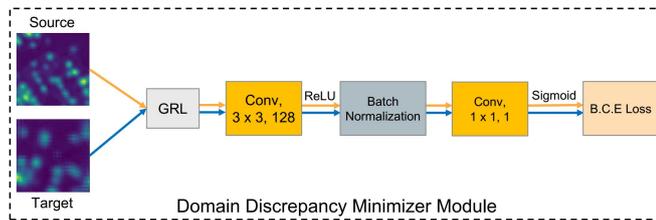
The counting pipeline, trained on the source dataset, produces the density map that has a structure similar to what is in the ground-truth maps. However, due to the domain shift, the distribution map produced for target data has visible artifacts as shown in Figure 7.

The learned features of the network are biased by the supervised training of the source domain. Thus the network does not recognize the features that it needs to localize the center of the building in the target domain. To learn a better output space distribution for the target domain, we must enforce the network to learn to output similar distribution for the target as for the source domain. For this, we perform adversarial alignment over the output distribution map using a domain discrepancy minimizer module as depicted in Figure 8. The process of adversarial alignment is described in Figure 9.

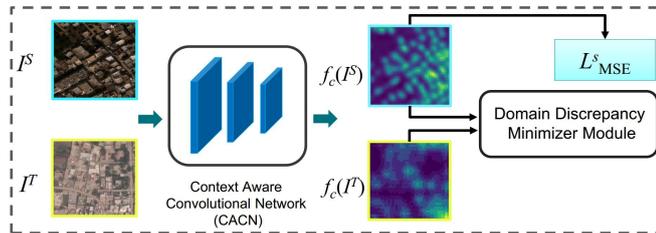
The output distribution of a network for two domains is considered to be consistent if we cannot discriminate between the two distributions. A domain discriminator can be used to learn to classify whether the generated distribution map is for the source image or target image. The more discriminative the two outputs are, the less consistent the two distributions would be and vice versa. Thus we can see that there is an inverse relationship between the consistency of outputs of the network and the ability of the discriminator to distinguish them. In adversarial domain adaptation, outputs for the two domains are aligned by forcing the network that produces the outputs to generate similar distributions. Since there is a relation between the discriminator and the density-maps-generating network, we can use the former to adjust the latter by using the loss gradients of the discriminator.



**Fig. 7** The distribution map predicted for the source images by source trained model has a visible structure consisting of multiple Gaussians. However, prediction over the target image patch results in a distribution map that is more blurred and lacks structure. Here,  $I^S$  is the source image patch,  $f_c(I^S)$  is the predicted distribution map of the source image patch from the source model,  $I^T$  is the target image patch,  $f_c(I^T)$  is the predicted distribution map of the target image patch from the source model,  $D^T$  is the ground-truth density map of the target image patch.



**Fig. 8** Domain Discrepancy Minimizer Module: A discriminator and GRL are added after the network outputs for adversarial domain adaptation.



**Fig. 9** This figure depicts distribution map alignment using adversarial learning.  $I^S$  is the source image patch.  $I^T$  is the target image patch.  $f_c(\cdot)$  represents density maps of their respective images.  $\mathcal{L}_{MSE}^s$  is the M.S.E loss between generated and ground truth density maps of the source image.

Ganin et al. [4] showed this the above-mentioned objective could be achieved by using a Gradient Reversal Layer (GRL) with a discriminator to learn domain invariant features. In this setting, the discriminator tries to distinguish the domains using a standard cross-entropy loss function. And as the training continues, we adjust the density-maps-generating network such that the discriminator is not able to distinguish the domains. This is where gradient reversal comes in. Gradients coming from the discriminator are reversed before being propagated to the density-maps-generating network. Thus, in effect, we are exploiting the reverse relation between the two where the gradients-updating-discriminator in one direction are updating the density-maps-generating-network in opposite direction. Thus while the discriminator is trying to distinguish the domains, the density-maps-

generating network is now trying to generate outputs such that the discriminator is unable to distinguish the domains. As the training continues, the discriminator finds it more and more difficult to distinguish the domains which means that the density maps generated are more and more consistent. Thus, we can align the outputs of the network for the two domains.

In our case, the discriminator is a small network containing a couple of convolution layers and a batch norm layer as shown in Figure 8. The discriminator minimizes the following standard classification loss function of Binary Cross Entropy:

$$\mathcal{L}_{B.C.E} = - \sum_{i=1}^{N_c} y_i \log(\hat{y}_i), \quad (2)$$

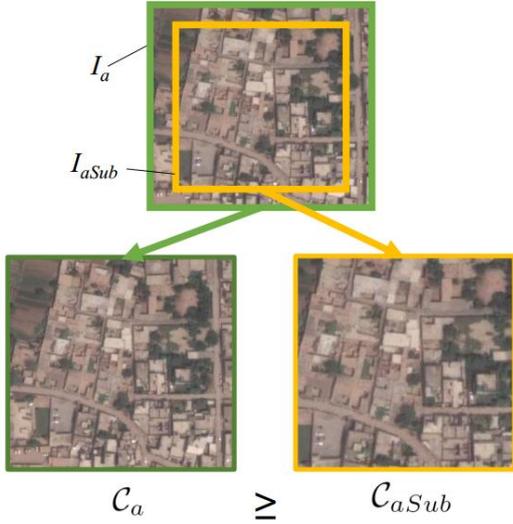
where  $\mathcal{L}_{B.C.E}$  is the Binary Cross Entropy Loss,  $y$  is the label vector and  $N_c$  is the number of classes which is two (source and target) in our case. The discriminator is applied after the output density maps and the GRL layer is used between the density maps outputs and the discriminator. It should also be noted that the discriminator and GRL are only added at training time to minimize the domain gap and are removed at test time. So, the total loss for the adversarial learning-based domain adaptation step is given as:

$$\mathcal{L}_{DMA} = \mathcal{L}_{MSE}^s + \alpha \mathcal{L}_{B.C.E}, \quad (3)$$

where  $\mathcal{L}_{DMA}$  is distribution map alignment loss,  $\alpha$  is the weighting factor and  $\mathcal{L}_{MSE}^s$  is the mean square error computed over the source domain (Equation 1).

#### 4.2.2 Counting Consistency within Image

To overcome the limitation of unlabeled target data, we design basic constraints that should be true for correct counting. In its basic form, any patch of the image cannot have more objects than the whole image. To accomplish this, we have employed the ranking loss. Ranking loss has been used



**Fig. 10** Within-Image consistency constraint ensures that the count of  $\mathcal{I}_a$  should be greater than or equal in value to the count of  $\mathcal{I}_{aSub}$ .

previously to constrain unsupervised deep learning-based methods [26]. Let  $C_a = C(f_c(\mathcal{I}_a))$  be predicted number of buildings in the image patch  $\mathcal{I}_a$  and the  $\mathcal{I}_{aSub}$  be the sub-patch (see Figure 10) extracted from  $\mathcal{I}_a$ , and  $C_{aSub} = C(f_c(\mathcal{I}_{aSub}))$  is the predicted number of buildings in the sub-patch. The within-image count consistency loss is given as:

$$\mathcal{L}_{WI}(\mathcal{I}_a, \mathcal{I}_{aSub}) = \max(0, -(\mathcal{C}_a - \mathcal{C}_{aSub}) + m), \quad (4)$$

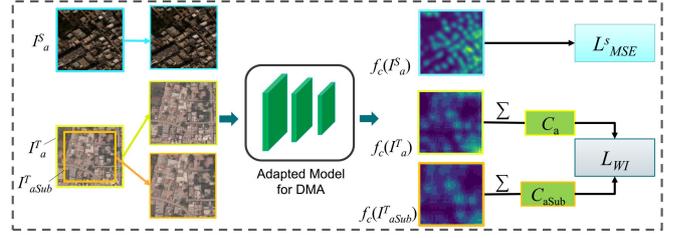
where  $\mathcal{L}_{WI}$  is the within-image counting consistency loss, margin  $m$  allows us to control the relaxation in the constraint. No loss is back-propagated in case of equal counts or when they are in the correct order as depicted in Equation 4. The total loss for this step is given as:

$$\mathcal{L}_{CWI} = \mathcal{L}_{DMA} + \lambda_1 \mathcal{L}_{WI}, \quad (5)$$

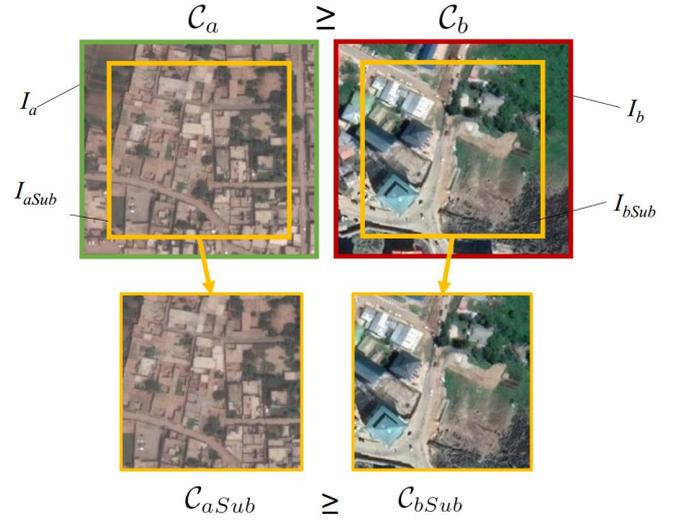
where  $\mathcal{L}_{CWI}$  is the total loss to compute within-image counting consistency,  $\mathcal{L}_{DMA}$  is distribution map alignment loss, and  $\lambda_1$  is the relative weight assigned. The whole process of learning employing counting consistency within images is shown in Figure 11.

#### 4.2.3 Counting Consistency Across Image

Counting consistency loss within the same images is not a powerful enough constraint, as indicated by [26] who used it for *warmup* task before using the supervised learning to predict the counting. Therefore, a much stronger across-image-counting-consistency constraint is proposed to help in domain adaptation. The constraint states if image patch  $\mathcal{I}_a$  has buildings substantially greater than the  $\mathcal{I}_b$ , any large-enough sub-patch of  $\mathcal{I}_a$  will also contain buildings greater than or equal to the number of buildings in an equally large-enough sub-patch of  $\mathcal{I}_b$  (see Figure 12). Let  $C_a = C(f_c(\mathcal{I}_a))$  and



**Fig. 11** This figure represents a learning framework employing image counting consistency constraint.  $I_a^S$  is the source image.  $I_a^T$  is the target image.  $I_a^TSub$  is the resized sub-image of  $I_a^T$ .  $f_c()$  represents density maps of their respective images.  $\mathcal{L}_{MSE}^S$  is the M.S.E loss between generated and ground truth density map of the source image.  $\mathcal{L}_{WI}$  is the within-image count consistency loss.



**Fig. 12** Counting consistency across images ensures that if the predicted count of  $\mathcal{I}_a$  is greater than or equal to the predicted count of  $\mathcal{I}_b$ , then the count of  $\mathcal{I}_{aSub}$  should also be greater than or equal to the count of  $\mathcal{I}_{bSub}$ .

$C_b = C(f_c(\mathcal{I}_b))$  be predicted number of buildings in the patches  $\mathcal{I}_a$  and  $\mathcal{I}_b$ , and  $C_{aSub}$  and  $C_{bSub}$  are the predicted number of buildings in the sub-patches respectively.

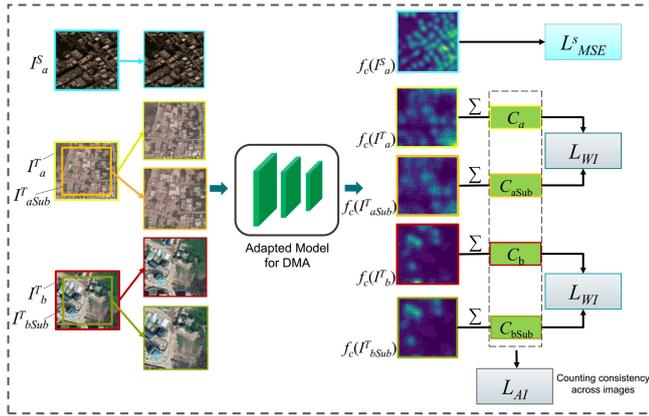
The constraint is represented as loss in the following equation,

$$\mathcal{L}_{AI}(\mathcal{I}_a, \mathcal{I}_b) = \begin{cases} \max(0, -(\mathcal{C}_{aSub} - \mathcal{C}_{bSub}) + m) & \text{if } \mathcal{C}_a \geq \mathcal{C}_b \\ \max(0, -(\mathcal{C}_{bSub} - \mathcal{C}_{aSub}) + m) & \text{otherwise} \end{cases} \quad (6)$$

where  $\mathcal{L}_{AI}$  is the across-image counting consistency loss,  $m$  is the margin. While implementing,  $\mathcal{I}_a$  and  $\mathcal{I}_b$  are chosen such that there is a larger than 5 difference in their count. To further keep this constraint true, the extracted sub-image is 80% of the original image. The total loss in the case of the across-image-consistency loss is given as:

$$\mathcal{L}_{CAI} = \mathcal{L}_{DMA} + \lambda_1 \mathcal{L}_{WI}(\mathcal{I}_a, \mathcal{I}_{aSub}) + \lambda_1 \mathcal{L}_{WI}(\mathcal{I}_b, \mathcal{I}_{bSub}) + \lambda_2 \mathcal{L}_{AI}(\mathcal{I}_a, \mathcal{I}_b), \quad (7)$$

where  $\mathcal{L}_{CAI}$  is the total loss to compute across-image counting consistency,  $\mathcal{L}_{DMA}$  is distribution map alignment loss,



**Fig. 13** This figure shows the learning framework using both within the image and across the image counting consistency constraints.  $I_a^S$  is the source image patch.  $I_a^T$ ,  $I_b^T$ ,  $I_{aSub}^T$ , and  $I_{bSub}^T$  are the target image patches and their respective resized sub-patches.  $f_c()$  represents density maps of their respective patches.  $\mathcal{L}_{MSE}^s$  is the M.S.E loss between generated and ground truth density map of the source image patch.  $\mathcal{L}_{WI}$  is the within-image count consistency loss.  $\mathcal{L}_{AI}$  is the across-image count consistency loss.

$\mathcal{L}_{WI}$  computes within-image counting consistency loss,  $\mathcal{L}_{AI}$  computes across-image counting consistency loss,  $\lambda_1$  and  $\lambda_2$  are the relative weights assigned. The whole process of maintaining counting consistency across images is depicted in Figure 13.

## 5 Experiments

### 5.1 Implementation details

Both source dataset and target datasets were divided into training, validation, and testing sets according to 60:20:20 ratios which made 2958 training images, 989 validation images, and 988 testing images. The image patches of both source and target datasets have a size of  $500 \times 500$  pixels each. As a pre-processing step, we performed histogram equalization on patches of both our source and target datasets such that they have the same contrast level. We trained the source model on 2958 training patches of xView dataset for 140 epochs using Adam as an optimizer and by keeping a learning rate of  $1e^{-4}$ .  $\mathcal{L}_{MSE}^s$  (Equation 1) is minimized between the predicted and ground truth density maps of the xView dataset. We kept the batch size to be equal to 26 patches. Validation was performed on 989 patches of the validation set of xView.

While adapting the source model for DMA, we have used the learning rate of  $1e^{-5}$  and  $\alpha$  was set to 0.1. During this adaptation process, the training set of xView images and the unlabeled training set of IML-DAC were utilized. In the experiments performed for within-image consistency only, the learning rate was kept at  $1e^{-5}$  and the total number of epochs was set to be 50. The relative weight  $\lambda$  was chosen

to be 45. During adaptation for both within-image consistency and across-the-images consistency, the learning rate and the number of epochs were  $1e^{-5}$  and 50, while  $\lambda_1$  and  $\lambda_2$  were taken to be 45 and 1 respectively. The batch size was again kept as 26 for these two adaptation experiments. The number of training images from both our source and target datasets was kept the same to be 2958 respectively. Hence a total of 5916 training image patches were utilized in the adaptation processes. Testing of these adapted models was performed on 988 patches of the testing set of the IML-DAC dataset and on 230 patches of the South Asian subset of xView dataset.

#### 5.1.1 Computation Cost:

Training is done on a single machine equipped with TitanX GPU having 12 GB memory. Source model training took approximately 12-13 hours to complete. Adaptation time for Distribution Map Alignment (DMA) took approximately 4 hours, within image Counting Consistency took approximately 4 hours, and the adaptation time for across-image counting consistency is approximately 6 hours. Testing on a single image takes (on average) 4678 ms.

### 5.2 Evaluation Metrics:

To evaluate our approach, we compute Mean Relative Error (MRE) :

$$MRE = \frac{1}{N} \sum_{i=1}^N \left( \frac{|C_{GT}(I_i) - C_{Pred}(I_i)| \times 100}{C_{GT}(I_i)} \right) \quad (8)$$

where  $N$  is the number of image patches in our testing set,  $C_{GT}(I_i)$  is the ground truth count of buildings in the  $i^{th}$  image and  $C_{Pred}(I_i)$  is the predicted count of buildings in the  $i^{th}$  image patch. MRE, also known as Mean Absolute Percentage Error (MAPE) [3, 44] is less susceptible to outliers in comparison to Mean Absolute Error (MAE) and the Root Mean Squared Error (RMSE) [3].

### 5.3 Hyper-parameter Selection

Histogram equalization is utilized as a pre-processing step to improve the contrast level of the images and is applied to each image separately. This preprocessing step improves the source trained model (that has not seen the target data) because it does not have to overcome the difference in contrast level in the two domains. In Figure 14, we have compared our patches of both xView and IML-DAC datasets, before and after applying histogram equalization. Below we present the results before and after performing histogram equalization. As shown in Table 5, the source-only trained model

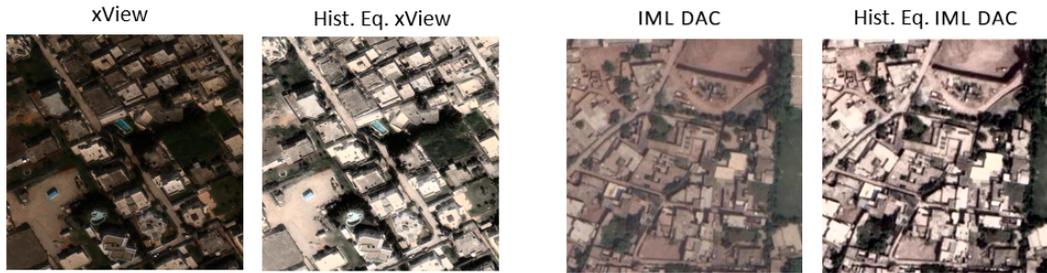


Fig. 14 Side-by-side comparison of patches, before and after histogram equalization, of IML-DAC and xView datasets.

Table 1 This table demonstrates the hyper-parameters selection of our different experiments.

Experiments	Parameters	MRE	$\omega$
$\mathcal{L}_{DMA}$	$\alpha = 0.05$	32.91	118
	$\alpha = 0.1$	<b>32.06</b>	<b>140</b>
	$\alpha = 0.15$	32.65	106
$\mathcal{L}_{CWI}$	$\alpha = 0.1, \lambda_1 = 35$	31.25	890
	$\alpha = 0.1, \lambda_1 = 40$	31.98	815
	$\alpha = 0.1, \lambda_1 = 45$	<b>30.16</b>	<b>1116</b>
	$\alpha = 0.1, \lambda_1 = 50$	30.94	1077
$\mathcal{L}_{CAI}$	$\alpha = 0.1, \lambda_1 = 45, \lambda_2 = 1/22$	29.77	1098
	$\alpha = 0.1, \lambda_1 = 45, \lambda_2 = 1/24$	28.86	1246
	$\alpha = 0.1, \lambda_1 = 45, \lambda_2 = 1/26$	<b>27.89</b>	<b>1464</b>
	$\alpha = 0.1, \lambda_1 = 45, \lambda_2 = 1/28$	28.37	1321

performed 22.9% better on the target dataset (IML-DAC) when histogram equalization was applied on patches.

Optimal hyper-parameters are selected by training on a small part of the training dataset in the target domain (IML-DAC) and testing over the full training dataset. The hyper-parameters were selected according to  $\omega$  which represents the number of image patches of the training set of our target data (IML-DAC) which followed within-image counting consistency and the MRE being computed on its testing set. Note that MRE is not used in choosing the parameters due to the assumption of the unavailability of ground truth density maps of the target training dataset. However, the correlation between the last two columns indicates the effectiveness of using with-in consistency loss for hyper-parameter selection. We start with  $\mathcal{L}_{DMA}$  (Eq. 3) and iterate over different values of  $\alpha$ . For each value of  $\alpha$ , the model is trained over a small dataset. The setting that results in the smallest within-image consistency loss over the full training dataset is chosen as the optimal value. Similarly, we find optimal value for  $\lambda_1$ , by minimizing  $\mathcal{L}_{CWI}$  (Eq. 5) and keeping  $\alpha$  constant. For  $\lambda_2$  both the  $\alpha$  and  $\lambda_1$  are kept equal to optimal values picked in previous steps, as we minimize  $\mathcal{L}_{CAI}$  (Eq. 7). In Table 1, our different sets of experiments demonstrate the usefulness of our selected hyperparameters.

## 5.4 Experimental Results

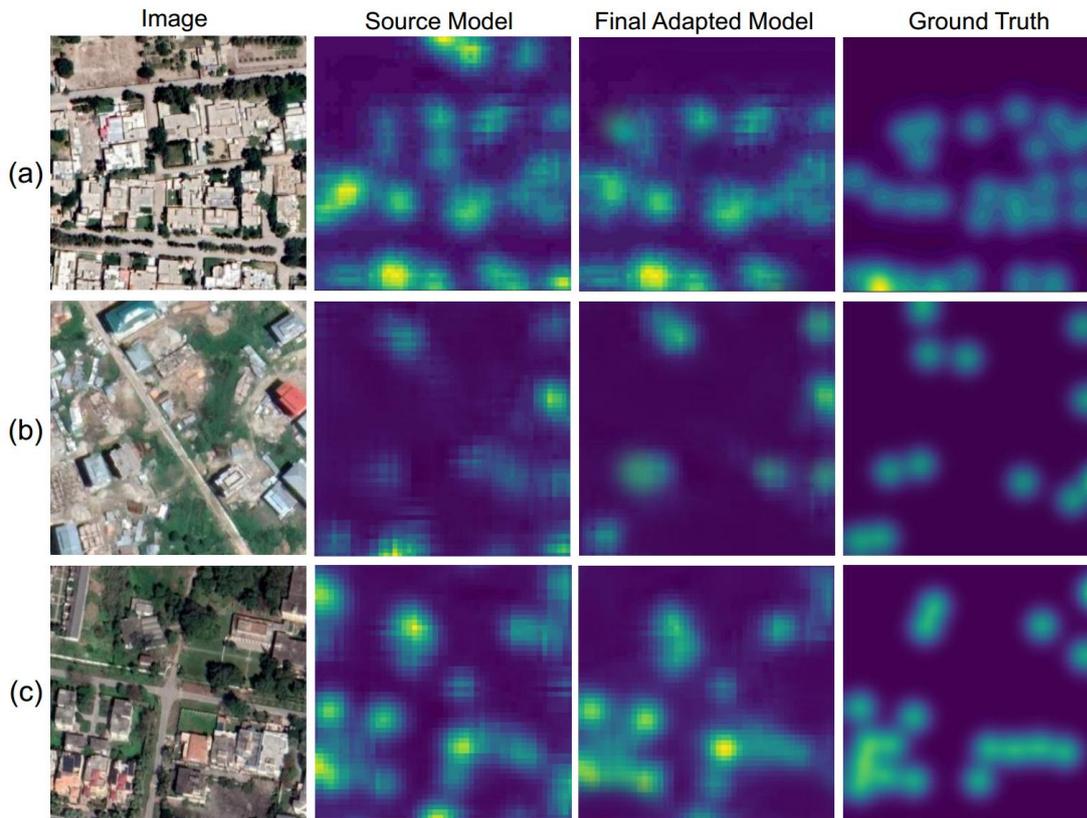
We evaluate our proposed approach on the testing set of IML-DAC which consists of 988 images and on a subset of xView consisting of South Asian countries.

**Component-wise analysis of the proposed approach:** A detailed quantitative comparison of the proposed approach and its components is given in Table 2 and Table 4. As indicated in Table 2, the source-trained model suffers a significant decline in performance when tested on the target (IML-DAC) dataset. We detail both, MRE over the target and the reduction in MRE with respect to when the only source-trained model is used, as different loss functions are introduced. Where density map alignment results in a decrease in Mean Relative Error (MRE), the significant improvement comes with the introduction of *Counting Consistency Constraints*, especially their combination. The final combination of all three losses results in the lowest counting error of 26.40%. A similar trend can also be seen while testing on south Asian regions of xView. Note that DMA also contains MSE loss.

Table 2 Comparison of mean relative error of all models on the target datasets. The first row is the model trained only on the source. Other rows represent adaptation by the indicated loss. Our final adapted model reduces error by approximately 7% on the IML-DAC dataset and approximately 20% on the South Asian subset of xView from the source trained model. MRE: lower is better. Reduction in Error: higher is better

Experiments	Target: IML-DAC		Target: xView (South Asian)	
	MRE	Reduction in Error	MRE	Reduction in Error
Source Trained Model	33.14%	-	48.24%	-
$\mathcal{L}_{DMA}$	31.97%	1.17%	37.40%	10.84%
$\mathcal{L}_{CWI}$	29.45%	3.69%	33.30%	14.94%
$\mathcal{L}_{CAI}$	<b>26.40%</b>	<b>6.74%</b>	<b>28.41%</b>	<b>19.83%</b>

To check the quality of density maps produced by our source and final adapted model, we have shown a detailed comparison in Figure 15. It is noticeable that as the model is adapted from the source to the target dataset, the predicted density maps improve in quality, i.e., they are better able to



**Fig. 15** From Left to Right: Input Image, distribution map predicted by source only trained model, distribution map predicted by the adapted model and Ground-truth. After adaptation, the predicted density map captures the localization information of the buildings much better than the ones produced by source only model.

locate and count the buildings. It is also worth noting that the density maps also become sharper as we adapt the model from the source to our target dataset. Moreover, in Figure 16 we observe the improvement in building counting for the satellite images using our proposed approaches.

**Comparison with related works:** Since, to the best of our knowledge, we are the first one to address the problem of building counting across the regions, we have compared our methods with previous methods which have addressed domain adaptation in crowd counting. The method of [7] minimized MMD loss between the source and target density maps, generated from crowded images, in a semi-supervised manner using three settings of the few-shot learning. In these three settings, 1, 5, and 10 labeled images respectively from the target domain were utilized while minimizing the MMD loss. We, however, have implemented their method in an unsupervised manner without using any labeled image of the target dataset. The work of [21] also proposes to address domain adaptation in crowd counting by utilizing ranking and adversarial loss to adapt the target dataset to cater to different density distributions and various scales.

**Table 3** In this table we compare our method with two of the existing works which deal with domain adaptation, but in crowd counting. Our final adapted model outperforms these models when tested on both the target datasets.

Experiments	Target: IML-DAC		Target: xView (South Asian)	
	MRE	Reduction in Error	MRE	Reduction in Error
Source Trained Model	33.14%	-	48.24%	-
MMD [7]	29.23%	3.91%	34.38%	13.86%
CODA [21]	28.17%	4.97%	30.67%	17.57%
<b>Ours</b>	<b>26.40%</b>	<b>6.74%</b>	<b>28.41%</b>	<b>19.83%</b>

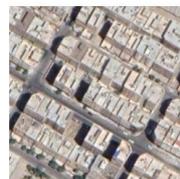
As demonstrated in Table 3, our proposed methodology outperforms both of these methods in terms of a higher reduction in error while testing on an unseen target domain.

**Results for different building count ranges:** To analyze the model’s behavior as the density of the buildings changes, we report results on different ranges in the number of buildings. For this purpose, we have segregated our testing data (IML-DAC) into three divisions: (i) images containing less than 31 buildings, (ii) images containing 31 to 60 buildings, (iii) images containing more than 60 buildings, and com-



Source Trained Model	65	56	18	56	12
$\mathcal{L}_{DMA}$	75	60	20	66	13
$\mathcal{L}_{CWI}$	76	78	25	70	16
$\mathcal{L}_{CAI}$	<b>94</b>	<b>100</b>	<b>31</b>	<b>86</b>	<b>22</b>
Ground Truth	96	97	29	85	23

**Fig. 16** Qualitative Results. Each column shows the improvement in building counting for the satellite image shown at the top of the column.

	Less than 31 Buildings		31 – 60 Buildings		More than 60 Buildings	
(a)	 GT = 16, Pred = 16	 GT = 29, Pred = 31	 GT = 33, Pred = 32	 GT = 53, Pred = 51	 GT = 85, Pred = 86	 GT = 97, Pred = 100
(b)	 GT = 20, Pred = 33	 GT = 20, Pred = 32	 GT = 60, Pred = 16	 GT = 41, Pred = 69	 GT = 96, Pred = 58	 GT = 64, Pred = 20

**Fig. 17** Comparison of predicted and ground truth counts using our final adapted model on images containing a different range of buildings. (a) depicts successfully predicted counts and (b) shows images where our model has failed to predict precise counts from images.

**Table 4** Comparison of mean relative error of all models across different ranges of buildings of target dataset (IML-DAC).

Building Ranges	Source Trained Model	$\mathcal{L}_{DMA}$	$\mathcal{L}_{CWI}$	$\mathcal{L}_{CAI}$
	MRE	MRE	MRE	MRE
Less than 31	30.65 %	29.15 %	26.57 %	<b>23.36 %</b>
31 - 60	51.03 %	52.83 %	51.13 %	<b>50.07 %</b>
More than 60	55.04 %	55.98 %	53.52 %	<b>50.78 %</b>

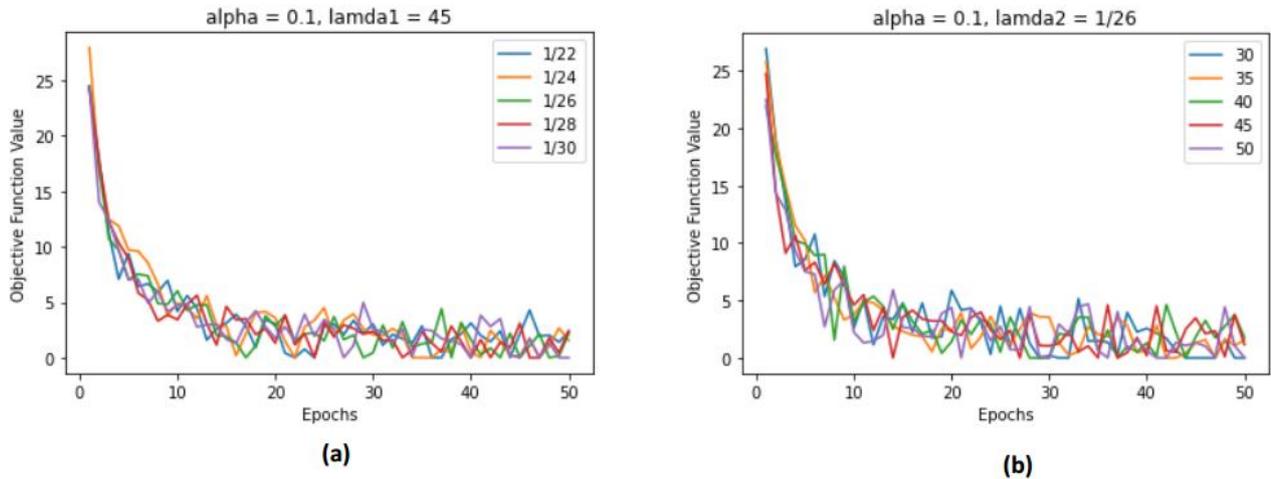
pared the results of different models in Table 4. We can observe that the count is more erroneous as we move to images with a high building count. In Figure 17, we present the qualitative results on some of the images from these three ranges of buildings. Figure 17(a) shows some images where count prediction is accurate, whereas Figure 17(b) shows cases where predicted counts are inconsistent from the ground truths.

## 5.5 Ablation

**Table 5** Comparison of M.R.E of source only model when trained on source dataset and tested on target dataset, with and without histogram equalization being part of preprocessing step. With histogram equalization generalization of the model improves.

Experiments	With Histogram Equalization	Without Histogram Equalization
	MRE	MRE
Source Trained Model Tested on IML-DAC	<b>33.14%</b>	42.98%

Histogram equalization is utilized to improve the contrast level of the images and is applied to each image separately. Figure 14, qualitatively show the effect of this step. This preprocessing step improves the source trained model (that has not seen the target data) because it does not have to overcome the difference in contrast level in the two do-



**Fig. 18** Convergence of performance during adaptation (a)  $\lambda_1$  is fixed to 45 and  $\lambda_2$  is being varied. (b)  $\lambda_2$  is fixed to 1/26 and  $\lambda_1$  is being varied.

mains. To show its effectiveness, we compared when the source-only trained model was trained on the source dataset and tested on the target dataset, without having this preprocessing step and when it is included. As shown in Table 2, the histogram equalization source-only trained model performed 22.9% better on the target dataset (IML-DAC).

The Figure 18 illustrates the convergence of our objective function. In these adaptation experiments shown in (a), we fixed alpha at 0.1, lambda1 ( $\lambda_1$ ) at 45, and varied lambda2 ( $\lambda_2$ ) from 1/22 to 1/30. The figure in (b) shows alpha fixed at 0.1, lambda2 ( $\lambda_2$ ) at 1/26 while lambda1 ( $\lambda_1$ ) is varied from 30 to 55. The given plots show the effect of these hyperparameters on performance convergence. *For all these values we see convergence, indicating that hyperparameters are not too sensitive in these ranges.*

## 6 Limitations & Future Directions

The current experiments on the target dataset were restricted to regions from South Asian countries only. The images covered not all but a few cities of these countries. For the future, a much larger dataset needs to be tagged and presented as standard for such studies, with special consideration to make it diverse and inclusive. In future work, we intend to include multi-task learning that exploits information such as the presence of roads or parks to improve the domain alignment and explainability component.

## 7 Conclusion

In this paper, we have addressed the challenging problem of cross-region building counting. We propose two counting consistency constraints to help direct the domain adap-

tation for the counting problem over the unlabeled target dataset. Exploiting the structure that should be there in the density map, we use adversarial learning to align the features across domains. Furthermore, we have introduced a large-scale dataset based on satellite imagery consisting of regions belonging to various South Asian regions to validate our domain adaptation methodology. The quantitative results prove that adapting the source trained model using our approach of count consistency and output space adaptation can predict counts from the target dataset quite accurately. Our proposed approach acts as a benchmark in this setting as it does not require any labeled images from the target dataset, making the whole process of building counting computationally efficient and labor-saving. We reported an improvement of approximately 7% on the IML-DAC dataset, and approximately 20% on the South Asian subset of xView over the model trained on the source dataset only. The huge improvement South Asian subset of xView could be due to the reason that both developed and developing regions in xView are captured at the same resolution and with the same sensors. On the other hand, comparatively less improvement on the IML-DAC dataset could be due to the large domain shift and demonstrate that our dataset is more challenging.

**Conflict of Interest:** We wish to confirm that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome. We confirm that the manuscript has been read and approved by all named authors and that there are no other persons who satisfied the criteria for authorship but are not listed. We further confirm that the order of authors listed in the manuscript has been approved by all of us. We confirm that we have given due consideration to the protection of intellectual property as-

sociated with this work and that there are no impediments to publication, including the timing of publication, with respect to intellectual property. We confirm that we have followed the regulations of our institutions concerning intellectual property. We understand that the corresponding author is the sole contact for the Editorial process (including the Editorial Manager and direct communications with the office). He/she is responsible for communicating with the other authors about progress, submissions of revisions, and final approval of proofs. We confirm that we have provided a current, correct email address which is accessible by the corresponding author.

**Data Availability:** The target dataset IML-DAC that supports the findings of our methodology can be accessed at: <https://github.com/intelligentMachines-ITU/domain-Adaptive-Building-Counting>.

## References

- Muhammad Usman Ali, Waqas Sultani, and Mohsen Ali. Destruction from sky: Weakly supervised approach for destruction detection in satellite imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 162:115–124, 2020.
- Bilel Benjdira, Yakoub Bazi, Anis Koubaa, and Kais Ouni. Unsupervised domain adaptation using generative adversarial networks for semantic segmentation of aerial images. *Remote Sensing*, 11(11):1369, 2019.
- Arnaud De Myttenaere, Boris Golden, Bénédicte Le Grand, and Fabrice Rossi. Mean absolute percentage error for regression models. *Neurocomputing*, 192:38–48, 2016.
- Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. PMLR, 2015.
- Ricardo Guerrero-Gómez-Olmedo, Beatriz Torre-Jiménez, Roberto López-Sastre, Saturnino Maldonado-Bascón, and Daniel Onoro-Rubio. Extremely overlapping vehicle counting. In *Iberian Conference on Pattern Recognition and Image Analysis*, pages 423–431. Springer, 2015.
- JT Harvey. Estimating census district populations from satellite imagery: Some approaches and limitations. *International journal of remote sensing*, 23(10):2071–2095, 2002.
- Mohammad Asiful Hossain, Mahesh Kumar Krishna Reddy, Kevin Cannons, Zhan Xu, and Yang Wang. Domain adaptation in crowd counting. In *2020 17th Conference on Computer and Robot Vision (CRV)*, pages 150–157. IEEE, 2020.
- Haopeng Hou, Yong Zhou, Jiaqi Zhao, Rui Yao, Ying Chen, Yi Zheng, and Abdulmoteleb El Saddik. Unsupervised cross-domain person re-identification with self-attention and joint-flexible optimization. *Image and Vision Computing*, 111:104191, 2021.
- Javed Iqbal and Mohsen Ali. Msl: Multi-level self-supervised learning for domain adaptation with spatially independent and semantically consistent labeling. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 1864–1873, 2020.
- Javed Iqbal and Mohsen Ali. Weakly-supervised domain adaptation for built-up region segmentation in aerial and satellite imagery. *ISPRS Journal of Photogrammetry and Remote Sensing*, 167:263–275, 2020.
- ISPRS. 2d semantic labelling contest. <https://www2.isprs.org/commissions/comm2/wg4/benchmark/semantic-labeling/>, accessed 24/12/2020,11:30 AM.
- Neal Jean, Marshall Burke, Michael Xie, W Matthew Davis, David B Lobell, and Stefano Ermon. Combining satellite imagery and machine learning to predict poverty. *Science*, 353(6301):790–794, 2016.
- Di Kang and Antoni Chan. Crowd counting by adaptively fusing predictions from an image pyramid. *arXiv preprint arXiv:1805.06115*, 2018.
- Nathan Korda, Balazs Szorenyi, and Shuai Li. Distributed clustering of linear bandits in peer to peer networks. In *International conference on machine learning*, pages 1301–1309. PMLR, 2016.
- Darius Lam, Richard Kuzma, Kevin McGee, Samuel Dooley, Michael Laielli, Matthew Klaric, Yaroslav Bulatov, and Brendan McCord. xview: Objects in context in overhead imagery. *arXiv preprint arXiv:1802.07856*, 2018.
- Guiying Li and Qihao Weng. Using landsat etm+ imagery to measure population density in indianapolis, indiana, usa. *Photogrammetric Engineering & Remote Sensing*, 71(8):947–958, 2005.
- He Li, Weihang Kong, and Shihui Zhang. Effective crowd counting using multi-resolution context and image quality assessment-guided training. *Computer Vision and Image Understanding*, 201:103065, 2020.
- 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 *2008 19th international conference on pattern recognition*, pages 1–4. IEEE, 2008.
- Shuai Li and Purushottam Kar. Context-aware bandits. *arXiv preprint arXiv:1510.03164*, 2015.
- Shuai Li, Alexandros Karatzoglou, and Claudio Gentile. Collaborative filtering bandits. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 539–548, 2016.
- Wang Li, Li Yongbo, and Xue Xiangyang. Coda: Counting objects via scale-aware adversarial density adaption. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 193–198. IEEE, 2019.
- Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan. Aggregating randomized clustering-promoting invariant projections for domain adaptation. *IEEE transactions on pattern analysis and machine intelligence*, 41(5):1027–1042, 2018.
- Jian Liang, Ran He, Zhenan Sun, and Tieniu Tan. Exploring uncertainty in pseudo-label guided unsupervised domain adaptation. *Pattern Recognition*, 96:106996, 2019.
- Wei Liu and Rongjun Qin. A multikernel domain adaptation method for unsupervised transfer learning on cross-source and cross-region remote sensing data classification. *IEEE Transactions on Geoscience and Remote Sensing*, 58(6):4279–4289, 2020.
- Weizhe Liu, Mathieu Salzmann, and Pascal Fua. Context-aware crowd counting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5099–5108, 2019.
- Xialei Liu, Joost Van De Weijer, and Andrew D Bagdanov. Leveraging unlabeled data for crowd counting by learning to rank. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7661–7669, 2018.
- Xu Liu, Steven W Chen, Shreyas Aditya, Nivedha Sivakumar, Sandeep Dcunha, Chao Qu, Camillo J Taylor, Jnaneshwar Das, and Vijay Kumar. Robust fruit counting: Combining deep learning, tracking, and structure from motion. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1045–1052. IEEE, 2018.
- Kanak Mahadik, Qingyun Wu, Shuai Li, and Amit Sabne. Fast distributed bandits for online recommendation systems. In *Proceedings of the 34th ACM international conference on supercomputing*, pages 1–13, 2020.
- Mark Marsden, Kevin McGuinness, Suzanne Little, Ciara E Keogh, and Noel E O’Connor. People, penguins and petri dishes:

- Adapting object counting models to new visual domains and object types without forgetting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8070–8079, 2018.
30. Azadeh Sadat Mozafari and Mansour Jamzad. Cluster-based adaptive svm: A latent subdomains discovery method for domain adaptation problems. *Computer Vision and Image Understanding*, 162:116–134, 2017.
  31. Daniel Onoro-Rubio and Roberto J López-Sastre. Towards perspective-free object counting with deep learning. In *European Conference on Computer Vision*, pages 615–629. Springer, 2016.
  32. Fang Qiu, Kevin L Woller, and Ronald Briggs. Modeling urban population growth from remotely sensed imagery and tiger gis road data. *Photogrammetric Engineering & Remote Sensing*, 69(9):1031–1042, 2003.
  33. Vincent Rabaud and Serge Belongie. Counting crowded moving objects. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 1, pages 705–711. IEEE, 2006.
  34. Maryam Rahnemoonfar and Clay Sheppard. Deep count: fruit counting based on deep simulated learning. *Sensors*, 17(4):905, 2017.
  35. Parvaneh Saeedi and Harold Zwick. Automatic building detection in aerial and satellite images. In *2008 10th International Conference on Control, Automation, Robotics and Vision*, pages 623–629. IEEE, 2008.
  36. Deepak Babu Sam, Shiv Surya, and R Venkatesh Babu. Switching convolutional neural network for crowd counting. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4031–4039. IEEE, 2017.
  37. Sanja Scepovic, Sagar Joglekar, Stephen Law, and Daniele Quercia. Jane jacobs in the sky: Predicting urban vitality with open satellite data. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW1):1–25, 2021.
  38. Anza Shakeel, Waqas Sultani, and Mohsen Ali. Deep built-structure counting in satellite imagery using attention based re-weighting. *ISPRS Journal of Photogrammetry and Remote Sensing*, 151:313–321, 2019.
  39. Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum self-paced learning for cross-domain object detection. *Computer Vision and Image Understanding*, 204:103166, 2021.
  40. Hang Su, Shaogang Gong, and Xiatian Zhu. Multi-perspective cross-class domain adaptation for open logo detection. *Computer Vision and Image Understanding*, 204:103156, 2021.
  41. M Naseer Subhani and Mohsen Ali. Learning from scale-invariant examples for domain adaptation in semantic segmentation. *arXiv preprint arXiv:2007.14449*, 2020.
  42. Waqas Sultani, Chen Chen, and Mubarak Shah. Real-world anomaly detection in surveillance videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6479–6488, 2018.
  43. Onur Tasar, SL Happy, Yuliya Tarabalka, and Pierre Alliez. Colormapgan: Unsupervised domain adaptation for semantic segmentation using color mapping generative adversarial networks. *IEEE Transactions on Geoscience and Remote Sensing*, 58(10):7178–7193, 2020.
  44. Chris Tofallis. A better measure of relative prediction accuracy for model selection and model estimation. *Journal of the Operational Research Society*, 66(8):1352–1362, 2015.
  45. Peter Tu, Thomas Sebastian, Gianfranco Doretto, Nils Krahnstoeber, Jens Rittscher, and Ting Yu. Unified crowd segmentation. In *European conference on computer vision*, pages 691–704. Springer, 2008.
  46. Adam Van Etten, Dave Lindenbaum, and Todd M Bacastow. Spacenet: A remote sensing dataset and challenge series. *arXiv preprint arXiv:1807.01232*, 2018.
  47. Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. Learning fine-grained image similarity with deep ranking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1386–1393, 2014.
  48. Litaow Wang, Shixin Wang, Yi Zhou, Wenliang Liu, Yanfang Hou, Jinfeng Zhu, and Futao Wang. Mapping population density in china between 1990 and 2010 using remote sensing. *Remote sensing of environment*, 210:269–281, 2018.
  49. Bo Wu and Ramakant Nevatia. Detection of multiple, partially occluded humans in a single image by bayesian combination of edgelet part detectors. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 1, pages 90–97. IEEE, 2005.
  50. Gui-Song Xia, Xiang Bai, Jian Ding, Zhen Zhu, Serge Belongie, Jiebo Luo, Mihai Datcu, Marcello Pelillo, and Liangpei Zhang. Dota: A large-scale dataset for object detection in aerial images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3974–3983, 2018.
  51. Gui-Song Xia, Jin Huang, Nan Xue, Qikai Lu, and Xiaoxiang Zhu. Geosay: A geometric saliency for extracting buildings in remote sensing images. *Computer Vision and Image Understanding*, 186:37–47, 2019.
  52. Feng Xiong, Xingjian Shi, and Dit-Yan Yeung. Spatiotemporal modeling for crowd counting in videos. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5151–5159, 2017.
  53. Christopher Yeh, Anthony Perez, Anne Driscoll, George Azzari, Zhongyi Tang, David Lobell, Stefano Ermon, and Marshall Burke. Using publicly available satellite imagery and deep learning to understand economic well-being in africa. *Nature communications*, 11(1):1–11, 2020.
  54. Laura Zabawa, Anna Kicherer, Lasse Klingbeil, Reinhard Töpfer, Heiner Kuhlmann, and Ribana Roscher. Counting of grapevine berries in images via semantic segmentation using convolutional neural networks. *ISPRS Journal of Photogrammetry and Remote Sensing*, 164:73–83, 2020.
  55. Jun Zhang, Jiao Liu, Bin Pan, and Zhenwei Shi. Domain adaptation based on correlation subspace dynamic distribution alignment for remote sensing image scene classification. *IEEE Transactions on Geoscience and Remote Sensing*, 58(11):7920–7930, 2020.
  56. Shanghang Zhang, Guanhang Wu, Joao P Costeira, and José MF Moura. Fcn-rlstm: Deep spatio-temporal neural networks for vehicle counting in city cameras. In *Proceedings of the IEEE international conference on computer vision*, pages 3667–3676, 2017.
  57. Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. Single-image crowd counting via multi-column convolutional neural network. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 589–597, 2016.
  58. Juepeng Zheng, Haohuan Fu, Weijia Li, Wenzhao Wu, Yi Zhao, Runmin Dong, and Le Yu. Cross-regional oil palm tree counting and detection via a multi-level attention domain adaptation network. *ISPRS Journal of Photogrammetry and Remote Sensing*, 167:154–177, 2020.
  59. Qiang Zhou, Shirui Wang, et al. Cluster adaptation networks for unsupervised domain adaptation. *Image and Vision Computing*, 108:104137, 2021.