

Bi-level Alignment for Cross-Domain Crowd Counting

Shenjian Gong¹, Shanshan Zhang^{1,*}, Jian Yang¹, Dengxin Dai², and Bernt Schiele²

¹PCA Lab, Key Lab of Intelligent Perception and Systems for High-Dimensional Information of Ministry of Education, and Jiangsu Key Lab of Image and Video Understanding for Social Security, School of Computer Science and Engineering, Nanjing University of Science and Technology

²MPI Informatics

{shenjiangong, shanshan.zhang, csjyang}@njust.edu.cn
{ddai, schiele}@mpi-inf.mpg.de

Abstract

Recently, crowd density estimation has received increasing attention. The main challenge for this task is to achieve high-quality manual annotations on a large amount of training data. To avoid reliance on such annotations, previous works apply unsupervised domain adaptation (UDA) techniques by transferring knowledge learned from easily accessible synthetic data to real-world datasets. However, current state-of-the-art methods either rely on external data for training an auxiliary task or apply an expensive coarse-to-fine estimation. In this work, we aim to develop a new adversarial learning based method, which is simple and efficient to apply. To reduce the domain gap between the synthetic and real data, we design a bi-level alignment framework (BLA) consisting of (1) task-driven data alignment and (2) fine-grained feature alignment. In contrast to previous domain augmentation methods, we introduce AutoML to search for an optimal transform on source, which well serves for the downstream task. On the other hand, we do fine-grained alignment for foreground and background separately to alleviate the alignment difficulty. We evaluate our approach on five real-world crowd counting benchmarks, where we outperform existing approaches by a large margin. Also, our approach is simple, easy to implement and efficient to apply. The code is publicly available at <https://github.com/Yankeegsj/BLA>.

1. Introduction

Crowd counting aims to estimate the number of persons in crowded scenes. This task has gained a lot of attention [10, 17] as it is useful for video surveillance, traffic control and human behavior analysis. Especially under pandemics such as COVID-19, it can be used to monitor and regulate

*Corresponding author

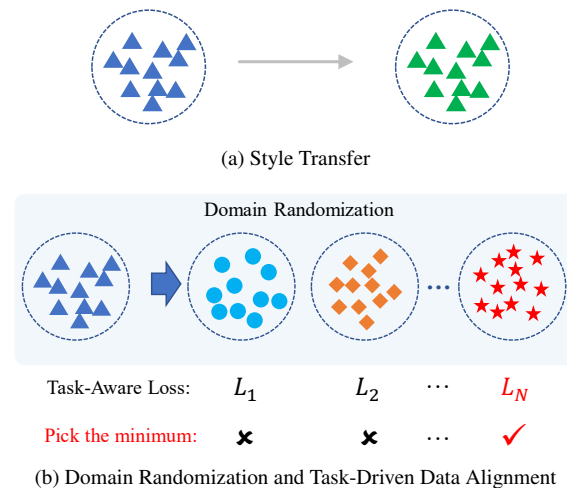


Figure 1. Comparison of three different ways for source domain augmentation. (a). Style transfer translates images from source to a target-like domain based on target style priors, but the translation is usually limited to color changes, and blind to the task objective. (b). Domain randomization augments the source domain randomly in a more diverse manner (colors, scales, etc.) but without any priors from target; our proposed task-driven data alignment is more similar to domain randomization; but instead of random selection, we pick the most suitable augmentation based on the task objective, which enables a more dynamic and robust model to the target domain.

the flow of people for safety reasons.

In recent years, many methods have been proposed for crowd counting. Most state-of-the-art approaches rely on ground truth density maps for a large number of training images, with each human head marked. However, it is extremely expensive to annotate so many human heads and for high-density regions such manual labels can be noisy. To reduce annotations costs, a large-scale synthetic dataset

GTA5 Crowd Counting (GCC) [24] was created, serving as a well-established training set with automatic annotations. However, models trained on the synthetic dataset perform poorly due to the large domain gap between synthetic and real-world images. Thus, it is necessary to investigate how to adapt the models trained on the synthetic domain to the real domain, without requiring annotations on the latter, i.e. via unsupervised domain adaptation (UDA).

There are a few UDA methods proposed for crowd counting. For instance, SE Cycle-GAN [24] translates synthetic images to the real domain with improved Cycle-GAN and then trains purely on the translated images; Gaussian Process-based iterative learning (GP) [22] generates pseudo-labels on the target images via a Gaussian process to allow for supervised training on the target domain. More recently, better performance has been achieved by employing an adversarial framework to align features from both source and target domains [5, 7]. However, FSC [7] introduces an auxiliary task of semantic segmentation, relying on external labeled human body segmentation datasets for pre-training; FADA [5] performs a coarse-to-fine estimation, making the inference less efficient.

In this paper, we aim to develop a new adversarial learning based method, which is more effective and flexible. We investigate the key components to boost performance. Previous methods employed either domain randomization or style transfer for source domain augmentation, using no priors or the target style priors only. In contrast, our task-driven data alignment is able to control the domain augmentation based on both the target style priors and the task objective such that it is optimized for our crowd counting task on the given target domain. We show a comparison of three different source domain augmentation methods in Fig. 1. On the other hand, since the foreground and background regions differ significantly in semantics, we propose a fine-grained feature alignment to handle them separately.

To summarize, the contributions of our work are as follows: (1) For more effective and efficient synthetic to real adaptive crowd counting, we propose a novel adversarial learning based method, consisting of bi-level alignments: task-driven data alignment and fine-grained feature alignment. (2) To the best of our knowledge, it is the first UDA approach to search for the optimal source data transform based on the downstream task performance on the target domain. (3) Experimental results on various real datasets show that our method achieves state-of-the-art results for synthetic-to-real domain adaptation; also, our method is simple and efficient to apply.

2. Related Works

Since we solve the problem of domain adaptive crowd counting, we first review recent works; our major contribution is a novel domain augmentation method via AutoML,

so we also discuss related works in the above two areas.

2.1. Domain Adaptive Crowd Counting

There are two groups of domain adaptation works in crowd counting: real-to-real and synthetic-to-real. Real-to-real adaptation aims to generalize models across real scenarios [8, 15], but since one real-world dataset is taken as the source domain, manual annotations are still needed. In contrast, synthetic-to-real adaptation fully avoids the requirement for manual annotations and thus is more interesting. In this work, we focus on synthetic-to-real adaptation. One direct way is to translate the labeled synthetic images to the style of the real images and then train on the translated images [4, 24], but it is limited by the performance of the translation method. Another intuitive way is to generate pseudo-labels on the target real images to allow for supervised learning on the target domain [4, 22]. More recently, the adversarial framework has been leveraged to achieve better performance via feature alignment between source and target domains [5, 7]. However, previous works are not efficient, requiring external training data or additional inference time. For instance, FSC [7] introduces an auxiliary task of semantic segmentation, relying on external labeled human body segmentation datasets for pre-training; FADA [5] performs a coarse-to-fine estimation, making the inference less efficient. In this work, we leverage adversarial training but aim to develop a simple yet effective method.

2.2. Domain Augmentation

Previously, there are two ways to augment the existing source domain: one is domain randomization, randomly changing the style of source images; the other is style transfer, translating the source images to the target style. Both of them preserve the contents of source images to allow for supervised learning.

Domain randomization is mostly used for domain generalization, which handles unknown target domains. The generated new samples with random styles are helpful to enhance the generalization ability of the trained model. One direct way is to generate various styles of images by style transfer CNN [29] or data augmentation [23]. On the other hand, AdaIN [12] has demonstrated that the mean and variance of convolutional feature maps can be used to represent the image style, making the domain randomization easier [19]. Some recent works modify styles in the frequency domain [11, 26] by randomizing the domain-variant components, which represent styles.

In contrast, style transfer is used for domain adaptation, where the target domain images are given to provide style guidance. The most common way is to generate source-to-target images through Cycle-GAN [4, 9, 24]. Similar to domain randomization, some works also do it in the frequency domain by replacing the style-related components with tar-

get ones, e.g. the amplitude of Fourier transform [28].

In this work, we propose a novel domain augmentation method named task-driven data alignment which is superior than domain randomization and style transfer. The major difference is that the augmentation is controlled by both the target style priors and by verifying the counting performance on a target-like domain. In this way, we are able to optimize augmentation for our crowd counting task on the given target domain.

2.3. Auto Machine Learning

Auto machine learning (AutoML) aims to free human practitioners and researchers from selecting the optimal values for each hyperparameter, such as learning rate, weight decay, and dropout [2], or designing well-performing network architectures [1]. Pioneers in this field develop optimization methods to guide the search process based on reinforcement learning (RL) [6], evolutionary algorithm (EA) [27] and Bayesian optimization [16]. These works are often impractical because of the required computational overhead. In contrast, a differentiable controller [18] converts the selection into a continuous hidden space optimization problem, allowing for an efficient search process performed by a gradient-based optimizer.

We apply a differentiable controller to search for several hyperparameters which represent styles. The transformed source images are then used for training, where the feature alignment becomes easier as the domain gap is reduced. In order to verify the searched hyperparameters, we construct a target-like domain in the feature space via AdaIN.

3. Bi-Level Alignment Counting (BLA)

In this section, we will introduce our bi-level alignment method for cross-domain crowd counting. Our core idea is to perform alignment between the source and target domains at both data-level and feature-level via two components namely task-driven data alignment and fine-grained feature alignment. The overall pipeline is depicted in Fig. 2 and detailed descriptions are provided in the following.

3.1. Problem Formulation

For UDA crowd counting, we have an annotated synthetic dataset $S = (x_S^i, y_S^i)_{i=1}^{N_S}$ as source and an unlabeled real-world dataset $T = (x_T^i)_{i=1}^{N_T}$ as target, where $x_S^i, x_T^i \in \mathbb{R}^{3 \times H \times W}$ denote an arbitrary image from the source and target domain, and $y_S^i \in \mathbb{R}^{H \times W}$ represents the ground truth density map in source. Our goal is to obtain a model that performs well on the target domain via reducing the large domain gap between the source and target.

3.2. Overview

As shown in Fig. 2, we propose a new UDA method for crowd counting based on adversarial learning. It consists

of feature extractor (\mathcal{F}), density estimator (\mathcal{E}), task-driven data alignment and local fine-grained discriminator (\mathcal{D}).

At training time, the source dataset S is first transformed to S^+ , with the same labels; a pair of images (x_{S^+}, x_T) from the augmented source and target domains are fed into \mathcal{F} , obtaining corresponding feature maps (F_S, F_T), F_S and $F_T \in \mathbb{R}^{C \times h \times w}$; \mathcal{D} performs feature alignment by passing reversed gradients to \mathcal{F} ; in the end, \mathcal{E} predicts density maps \tilde{y}_S based on F_S , supervised by y_S . At test time, the inference is rather simple: each target image x_T^i only goes through \mathcal{F} and \mathcal{E} to obtain the predicted density map \tilde{y}_T .

Following previous works, we employ VGG16 [20] as our feature extractor \mathcal{F} . For \mathcal{E} , we stack a series of convolution and deconvolution layers, inspired by [4].

The density estimation loss \mathcal{L}_E on a labeled source image can be defined as follows:

$$\mathcal{L}_E = \sum \|y_S - \tilde{y}_S\|^2. \quad (1)$$

3.3. Task-Driven Data Alignment

Typically, domain augmentation enhances the performance on the target domain using generated new samples via domain randomization or style transfer. In contrast to previous domain augmentation methods that are blind to the downstream task, our method searches for the most suitable augmentation based on both the target styles and the task performance on target via AutoML. In such a task-driven way, our method is expected to find a transform that better serves for the downstream task on the given target domain.

It has been shown in [17] that for the crowd counting task, images mainly differ in color, scale and perspective. Accordingly, in this paper, we define three basic transform units: RGB2Gray, scaling and perspective transform. Each transform is a combination of the above three transform units. As shown in Tab. 1, one transform is defined by five parameters, among which only three are searched for simplicity. It is notable the transform is not limited to the above three units and can be easily extended to different and also more types, with proper manual definitions.

Transform unit	Parameters	
	Split ratio	Attribute
RGB2Gray	pg*	-
Scaling	ps	Scale factor*
Perspective transform	ppt	Angle*

Table 1. Each transform consists of three different units, each represented by two parameters: one for split ratio and another for attribute. * marks those parameters we search while others are fixed. A full transform set is generated by iterating each parameter.

Given a transform, we split the whole source set into sev-

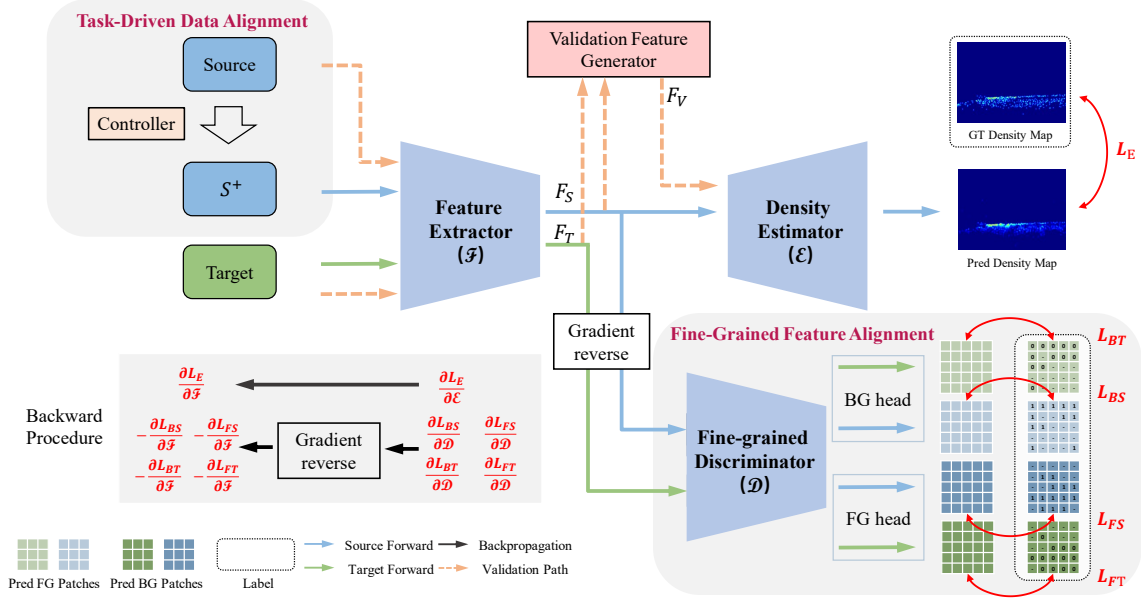


Figure 2. Overview of our proposed bi-level alignment framework (BLA), which mainly consists of four components: feature extractor (\mathcal{F}), density estimator (\mathcal{E}), task-driven data alignment and local fine-grained discriminator (\mathcal{D}). At training time, the source dataset S is transformed to S^+ with the optimal transform searched via task-driven data alignment (Alg. 1), during which the validation feature generator provides target-like features for candidate transform validation. Then the entire network is optimized based on S^+ and T using the training objective in Eq. 5. At test time, we simply feed a target image x_T^i to \mathcal{F} and \mathcal{E} to obtain the predicted density map \tilde{y}_T .

eral subsets via a transform tree, as shown in Fig. 3. At the 1st level, the whole source dataset is split into two subsets with a ratio of p_G , i.e. some images are converted to gray scale images (along path \mathbf{Y}), while others are kept the same (along path \mathbf{N}). At the following levels, each subset generated from the previous level goes through one split with a given split ratio and an attribute parameter. Given three transform units, we finally generate 8 (2^3) subsets.

In this paper, we use a differential controller to guide the direction of our search process. The search process is iterated via multiple rounds, each of which is described in Alg. 1. At each round, we first transform the source data given some transform candidates, and then obtain the reward of each transform via validation on a generated target-like set; after that, we learn the mapping function from transforms to corresponding rewards via training a differentiable controller; finally, we update the transform candidates based on the controller and goes to the next search round. In the following, we explain the above process in more detail.

Source Data Transformation Based on Candidate Transforms. At this stage, we first randomly initialize a transform set $\mathbb{D} = (d_k)_{k=1}^{N_D}$. Given an arbitrary transform d_k , we split the whole source set into several subsets via a transform tree, as illustrated in Fig. 3. In this way, we apply one transform on the source data and obtain a new mixed source dataset $S_k^+ = (x_{S_k^+}^i, y_{S_k^+}^i)_{i=1}^{N_{S^+}}$. Due to the large num-

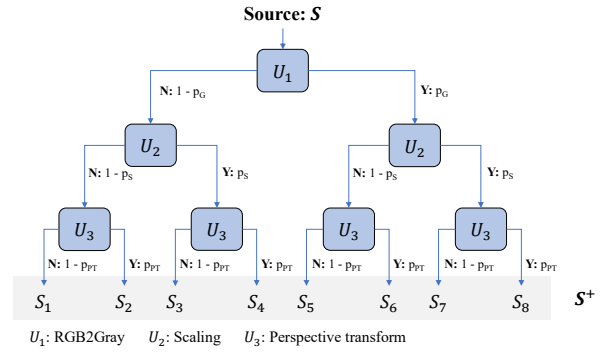


Figure 3. We split the whole source set into several subsets via a transform tree according to p_G , p_S and p_{PT} . “ $\mathbf{Y}: p_G$ ” refers to performing the transformation G with a ratio of p_G ; “ $\mathbf{N}: 1 - p_G$ ” means that no transformation is performed with a ratio of $1 - p_G$. Finally, we obtain the corresponding subsets as shown at the bottom row, which constitute S^+ .

ber of source images, we do not do standard data augmentation, but apply the transform on the original data to keep the same size of S_k^+ and S .

Candidate Transform Validation. Based on each new source dataset S_k^+ , we train the whole network as shown in Fig. 2 with the learning objective from Eq. 5. Now we need to evaluate the alignment quality of each S_k^+ . Ideally,

this should be done by measuring the counting performance on T . Unfortunately, we do not have labels for T . To address this problem, we propose a validation feature generator, which takes the features of a pair of source and target image features (F_S, F_T) from the feature extractor as input and generate a new feature F_V via AdaIN [12], which is a mixture of source contents and target style, namely a target-like image feature.

Specifically, we first compute the source style representation with channel-wise mean and standard deviation $\mu(F_S), \sigma(F_S) \in \mathbb{R}^C$ and the target style representation $\mu(F_T), \sigma(F_T) \in \mathbb{R}^C$. Then we replace the style of F_S with that of F_T and obtain F_V :

$$F_V^c = \mu(F_T^c) + \sigma(F_T^c) \cdot \left(\frac{F_S^c - \mu(F_S^c)}{\sigma(F_S^c)} \right), \quad (2)$$

where $c \in \{1, 2, 3, \dots, C\}$ is the channel index. After that, we feed F_V to the density estimator and get \tilde{y}_V . Because F_S and F_V share the same contents, we evaluate \tilde{y}_V based on y_S according to Eq. ???. In this way, we obtain the evaluated validation performance p_k as reward for transform d_k .

Candidate Transform Update. After obtaining the reward for each transform in \mathbb{D} , we then train a differentiable controller and let it learn the mapping function from a transform to its corresponding reward. The controller is of encoder-decoder structure. The encoder takes a transform as input, maps it to a hidden state, and predicts its performance as \tilde{p}_k . The decoder reconstructs the transform d_k as \tilde{d}_k from the hidden state. The loss function of our controller is defined as:

$$\mathcal{L}_C = \left\| d_k - \tilde{d}_k \right\|^2 + \left\| p_k - \tilde{p}_k \right\|^2. \quad (3)$$

Same with NAO [18], we then update the hidden state towards the gradient direction of improved performance and obtain a new transform set \mathbb{D}' , for better alignment. After several rounds, we choose the optimal transform from all validated transforms based on their rewards. Please refer to [18] for more details regarding the update procedure.

3.4. Fine-Grained Feature Alignment

To perform feature alignment, we employ adversarial learning via a discriminator and a gradient reverse layer. Inspired by the success of using segmentation as an auxiliary task for crowd counting [21], we propose a fine-grained discriminator, with two separated classification heads for foreground and background regions. To handle the unbalanced numbers of foreground and background pixels, the discriminator is applied on local patches instead of pixels.

Given the grid size $G = (g_h, g_w)$, we feed a pair of feature maps (F_S, F_T) to \mathcal{D} and obtain two pairs of patch-level discrimination maps: (O_{FS}, O_{BS}) for source and (O_{FT}, O_{BT}) for target, separating foreground

Algorithm 1 Pseudo code of one-round search procedure of data-level alignment

Input: Source and target domain training set S, T ; the pre-trained source only network $\hat{\mathcal{G}} = \{\hat{\theta}_{\mathcal{F}}, \hat{\theta}_{\mathcal{E}}\}$; the candidate transform set \mathbb{D} and the controller \mathcal{C}

```

1: for  $d_k \in \mathbb{D}$  do
2:   initialize  $\mathcal{G} = \{\theta_{\mathcal{F}}, \theta_{\mathcal{E}}, \theta_{\mathcal{D}}\}$  with  $\hat{\mathcal{G}}$ 
3:    $S_k^+ = \text{transform}(S, d_k)$  {As Fig. 3}
4:    $\mathcal{G} = \text{train}(\mathcal{G}, S_k^+, T)$ 
5:    $p_k = \text{validate}(\mathcal{G}, S, T)$  {Val performance as reward}
6: end for
7:  $\mathbb{P} = \{p_1, p_2, \dots\}$ 
8:  $\mathcal{C} = \text{train controller}(\mathcal{C}, \mathbb{D}, \mathbb{P})$  {As Eq. 3}
9:  $\mathbb{D} = \text{update}(\mathcal{C}, \mathbb{D})$  {Same with NAO [18]}
10: return  $\mathbb{D}, \mathbb{P}$ 

```

and background. Each map $O_{FS}, O_{FT}, O_{BS}, O_{BT} \in \mathbb{R}^{(H/g_h) \times (W/g_w)}$.

The segmentation masks (M_S, \tilde{M}_T) are obtained by thresholding the ground truth density maps. Please note that for the target we use pseudo density maps instead. Specifically, we apply a threshold of th on each patch, to threshold the sum of all its pixel values.

As shown in Fig. 2, we define local fine-grained discrimination losses of background and foreground $\mathcal{L}_B, \mathcal{L}_F$, as follows:

$$\begin{aligned}
\mathcal{L}_D &= \mathcal{L}_B + \mathcal{L}_F \\
\mathcal{L}_B &= \mathcal{L}_{BS} + \mathcal{L}_{BT} \\
&= \sum - (1 - M_S) \cdot \log(1 - O_{BS}) \\
&\quad + \sum - (1 - \tilde{M}_T) \cdot \log(O_{BT}), \quad (4) \\
\mathcal{L}_F &= \mathcal{L}_{FS} + \mathcal{L}_{FT} \\
&= \sum - M_S \cdot \log(1 - O_{FS}) \\
&\quad + \sum - \tilde{M}_T \cdot \log(O_{FT}).
\end{aligned}$$

We use the same back-propagation optimizing scheme with the gradient reverse layer [3] for adversarial learning.

3.5. Optimization

The optimization objective of the whole method is:

$$\mathcal{L} = \mathcal{L}_E + \lambda \mathcal{L}_D, \quad (5)$$

where λ is a weight factor to balance the task loss \mathcal{L}_E and the domain adaptation loss \mathcal{L}_D .

The whole network is optimized via two steps. For data alignment, we first optimize all the parameters in the network including the feature alignment component for each transform to obtain the corresponding reward, during the

search process in Alg. 1. After selecting the best transform, we retrain the whole network including the feature alignment component with the transformed S^+ .

4. Experiments

We first introduce the datasets, evaluation metrics and implementation details; then we provide comparisons with state-of-the-art methods, followed by analysis on data alignment; finally, we perform some ablation studies.

4.1. Datasets and Evaluation Metrics

To evaluate the proposed method, the experiments are conducted under adaptation scenarios from GCC [24] to five large-scale real-world datasets, i.e. ShanghaiTech Part A/B (SHA/S HB) [31], QNRF [14], UCF-CC-50 [13] and WorldExpo’10 [30] respectively. Statistics are listed in Tab. 2.

Dataset	Attribute	# Images	Cnt (Mean \pm Std)
GCC	Syn	15,211	501 \pm 718
SHA	Real	482	501 \pm 456
SHB	Real	716	123 \pm 94
QNRF	Real	1,535	815 \pm 1176
UCF-CC-50	Real	50	1,279 \pm 950
WorldExpo’10	Real	3,980	50 \pm 41

Table 2. Statistics of five real-world (Real) datasets and one synthetic (Syn) dataset GCC used for experiments.

Following previous works, we adopt Mean Absolute Error (MAE) and Mean Squared Error (MSE) as evaluation metrics. They are formulated as: $MAE = \frac{1}{N} \sum_{i=1}^N |\sum y_i - \sum \tilde{y}_i|$, and $MSE = \sqrt{\frac{1}{N} \sum_{i=1}^N |\sum y_i - \sum \tilde{y}_i|^2}$, where N is the number of test images; $\sum y_i$, $\sum \tilde{y}_i$ represent the ground truth and predicted number on the i -th image respectively.

4.2. Implementation Details

The architectures of feature extractor (\mathcal{F}), density estimator (\mathcal{E}), fine-grained discriminator (\mathcal{D}) and controller (\mathcal{C}) are listed in supplementary. We input 4 pairs of source and target images with a uniform size of 576×768 at each iteration. Following the previous work [4], we generate the ground truth density map using Gaussian kernel with a kernel size of 15×15 and a fixed standard deviation of 4. We set th , G in Eq. 4 and λ in Eq. 5 to 0.005, (16,16) and 1.0 respectively for all our experiments; both p_S and p_T are set to 0.5. And the gradient factor of the gradient reverse layer is set to 0.01. We also adopt a scene regularization strategy proposed by [24] to avoid negative knowledge transfer. We train the adversarial framework and controller with

Adam optimizer with default parameters, and their learning rates are initialized as 10^{-5} and 10^{-1} respectively. All experiments are conducted on a single NVIDIA RTX 2080TI GPU with 11GB of VRAM and our code is implemented with Pytorch.

4.3. Comparisons with State-of-the-Art

We compare our method BLA with previous published unsupervised domain adaptive crowd counting methods under the adaptation scenarios from synthetic GCC dataset to five different real-world datasets. All methods employ VGG16 [20] as backbone.

From the results in Tab. 3, we have the following observations: (1) Our proposed method outperforms all existing domain adaptation methods by a large margin across different datasets and on WorldExpo’10 we achieve comparable results with DACC. In particular, on SHA our proposed method achieves 99.3 MAE and 145.0 MSE, outperforming previous best results by 13.1 pp w.r.t. MAE and 31.9 pp w.r.t. MSE. (2) Our method is robust across various real target datasets, showing high adaptability. As shown in Tab. 2, although the density of these real-world datasets varies a lot, we perform the best on all target domains.

4.4. Analysis on Task-Driven Data Alignment

In order to understand how the data alignment behaves, we show the searched transforms for different target datasets in Tab. 4. It shows that our task-driven data alignment is quite interpretable, representing various domain gaps between source and different target domains. For instance, GCC contains highly-saturated color images while UCF-CC-50 and WorldExpo’10 contain lots of images with low saturations, so the RGB2Gray ratios on them are rather high (0.85 and 0.98) as gray scale images are of 0 saturation; in contrast, SHB is closer to GCC in terms of saturation, so the RGB2Gray ratio on SHB is rather low (0.16). Similarly, since UCF-CC-50 is denser than other datasets, its scaling factor is particularly smaller, such that denser regions with small-scale heads will be generated.

Moreover, we perform cross dataset validation by applying the searched transform on one dataset to another. As shown in Tab. 5, when testing on SHB, the transform searched on SHA underperforms that searched on SHB; and vice versa. These results indicate that each dataset requires its personalized transform to achieve the optimal performance. It is necessary to search for the most suitable transform on each dataset in an automatic way so as to avoid tedious manual designs.

Additionally, Fig. 4 shows some qualitative results on the SHA dataset. From Column 3, we can see that without adaptation, the model either fails to detect the presence of people in some areas (top row), or fails to get a correct estimate of the local density (middle and bottom rows). Our

Method	SHA		SHB		QNRF		UCF-CC-50		WorldExpo'10
	MAE	MSE	MAE	MSE	MAE	MSE	MAE	MSE	Avg. MAE
NoAdpt [32]	160.0	216.5	22.8	30.6	275.5	458.5	487.2	689.0	42.8
Cycle-GAN [32]	143.3	204.3	25.4	39.7	257.3	400.6	404.6	548.2	26.3
SE Cycle-GAN [24]	123.4	193.4	19.9	28.3	230.4	384.5	373.4	528.8	26.3
SE Cycle-GAN(JT) [25]	119.6	189.1	16.4	25.8	225.9	385.7	370.2	512.0	24.4
FSC [7]	129.3	187.6	16.9	24.7	221.2	390.2	-	-	-
FADA [5]	-	-	16.0	24.7	-	-	-	-	21.6
GP [22]	121.0	181.0	12.8	19.2	210.0	351.0	355.0	505.0	20.4
DACC [4]	112.4	176.9	13.1	19.4	211.7	357.9	-	-	17.4
BLA (ours)	99.3	145.0	11.9	18.9	198.9	316.1	346.8	480.0	17.9

Table 3. Comparison of our method with previous methods for synthetic-to-real adaptation. All methods employ VGG16 [20] as backbone.

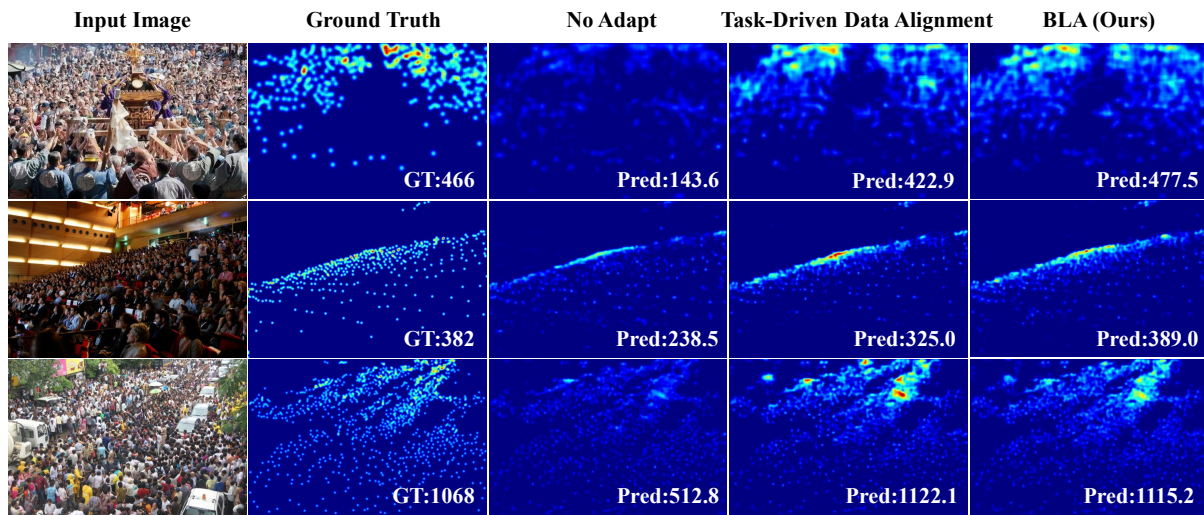


Figure 4. Qualitative adaptation results of images with different levels of density on SHA.

Dataset	Searched transform d		
	p_G	Scale factor	Angle
SHA	0.78	0.77	12°
SHB	0.16	0.42	3°
QNRF	0.67	0.41	25°
UCF-CC-50	0.85	0.23	4°
WorldExpo'10	0.98	0.57	17°

Table 4. Searched transforms vary across different target datasets.

task-driven data alignment helps to reduce the errors largely, indicating the domain gaps are significantly narrowed. After further adding our fine-grained feature alignment, our BLA method provides more accurate final counts.

Dataset	Searched d for SHA		Searched d for SHB	
	[0.78, 0.77, 12°]		[0.16, 0.42, 3°]	
	MAE	MSE	MAE	MSE
SHA	99.3	145.0	104.3	153.7
SHB	14.7	26.4	11.9	18.9

Table 5. Different datasets require different transformations for optimal results.

4.5. Ablation Studies

In this subsection, we conduct some ablation studies to analyze different components of our proposed BLA. All experiments are conducted under the GCC \rightarrow SHA adaptation due to its large variation in crowd density.

Effects of Two Levels of Alignment. We first analyze the effects of data alignment and feature alignment.

As shown in Tab. 6, the performance is significantly improved by ~ 25 pp w.r.t MAE (from 134.7 to 109.1) when task-driven data alignment is employed. On the other hand, we also observe a large improvement of ~ 14 pp w.r.t MAE (from 134.7 to 121.1) by replacing global feature alignment with fine-grained feature alignment. Moreover, we obtain a total gain of ~ 35 pp w.r.t MAE by adding both alignments. These results indicate the effects of two levels of alignment and the complementarity between them.

Task-Driven Data Alignment	Fine-Grained Feature Alignment	MAE	MSE
×	×	134.7	210.9
✓	×	109.1	153.8
×	✓	121.1	200.8
✓	✓	99.3	145.0

Table 6. Effects of two levels of alignment.

Effect of Task-Driven Data Alignment To evaluate the effectiveness of task-driven data alignment, we replace it with domain randomization and style transfer. From Tab. 7, we can see that our task-driven data alignment outperforms previous two data augmentation methods by a large margin.

Method	MAE	MSE
Style Transfer	119.4	194.6
Domain Randomization	110.0	164.2
Task-Driven Data Alignment	99.3	145.0

Table 7. Effect of Task-Driven Data Alignment.

Effect of Validation Feature Generator. To evaluate the alignment quality of each transform, we generate a target-like feature set for validation via AdaIN. In Fig. 5, we compare the counting performance on our generated validation set and the real target training set w.r.t. MAE, where the index indicates different combinations of transform parameters. We can see that the two curves go in a similar trend, i.e. the worst performance happens at index 0, the best performance happens at index 7, and there is fluctuation in between. This comparison demonstrates that our generated validation set is of high similarity to the real target set, allowing us to do effective validation without relying on target annotations. On the other hand, we observe the performance varies a lot along the choice of transform parameters, showing that different transforms highly affect the performance. Thus it is of great importance to search for an optimal transform for a given target set.

Impact of Grid Size in Fine-Grained Feature Alignment. Our fine-grained feature alignment strategy is conducted in a patch-wise style. We analyze how the grid size

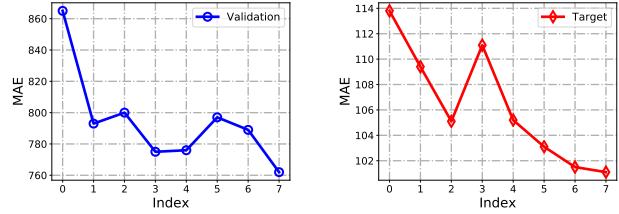


Figure 5. Comparison of validation performance on the generated set (left) and real target set (right) across different transforms w.r.t. MAE. The similar trend verifies the effect of our validation feature generator.

G affects the performance. As shown in Tab. 8, if the size G is too small or too large, there will be a data imbalance between the numbers of background and foreground patches, which results in poor feature alignment. Since the patch size of 16 performs the best, we use 16 as the default size in all experiments.

Grid size G	MAE	MSE
(2,2)	138.1	222.9
(4,4)	130.5	205.4
(8,8)	129.0	206.3
(16,16)	121.1	200.8
(32,32)	141.4	223.6

Table 8. Impact of local grid size G used in fine-grained feature alignment.

In the supplementary material, we provide more ablation studies on the impact of segmentation threshold, effect of additional style transfer from S^+ to T , effect of using more transformations, and a comparison to grid search.

5. Conclusion

In this work, we propose a bi-level alignment framework for synthetic-to-real UDA crowd counting. On one hand, we propose task-driven data alignment to search for a specific transform given the target set, which is applied on the source data to narrow down the domain gap at the data level. On the other hand, to alleviate the alignment difficulty on the entire image, we propose to perform fine-grained feature alignment on foreground and background patches separately. Extensive experiments on five real-world crowd counting benchmarks have demonstrated the effectiveness of our contributions.

Acknowledgements

This work was supported in part by the “111” Program B13022, Fundamental Research Funds for the Central Universities (No. 30920032201) and the National Natural Science Foundation of China (Grant No. 62172225).

References

- [1] Minghao Chen, Jianlong Fu, and Haibin Ling. One-shot neural ensemble architecture search by diversity-guided search space shrinking. In *CVPR*, pages 16530–16539, 2021. 3
- [2] Xuanyi Dong, Mingxing Tan, Adams Wei Yu, Daiyi Peng, Bogdan Gabrys, and Quoc V Le. Autohas: Differentiable hyper-parameter and architecture search. In *arxiv*, 2020. 3
- [3] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *Int. Conf. Mach. Learning*, pages 1180–1189, 2015. 5
- [4] Junyu Gao, Tao Han, Qi Wang, and Yuan Yuan. Domain-adaptive crowd counting via inter-domain features segregation and gaussian-prior reconstruction. In *arxiv*, 2019. 2, 3, 6, 7
- [5] Junyu Gao, Yuan Yuan, and Qi Wang. Feature-aware adaptation and density alignment for crowd counting in video surveillance. In *arxiv*, 2019. 2, 7
- [6] Xinyu Gong, Shiyu Chang, Yifan Jiang, and Zhangyang Wang. Autogan: Neural architecture search for generative adversarial networks. In *ICCV*, pages 3224–3234, 2019. 3
- [7] Tao Han, Junyu Gao, Yuan Yuan, and Qi Wang. Focus on semantic consistency for cross-domain crowd understanding. In *ICASSP*, pages 1848–1852, 2020. 2, 7
- [8] Yuhang He, Zhiheng Ma, Xing Wei, Xiaopeng Hong, Wei Ke, and Yihong Gong. Error-aware density isomorphism reconstruction for unsupervised cross-domain crowd counting. In *AAAI*, number 2, pages 1540–1548, 2021. 2
- [9] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. Cycada: Cycle-consistent adversarial domain adaptation. In *Int. Conf. Mach. Learning*, pages 1989–1998, 2018. 2
- [10] 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
- [11] Jiaying Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. Fsd: Frequency space domain randomization for domain generalization. In *CVPR*, pages 6891–6902, 2021. 2
- [12] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1501–1510, 2017. 2, 5
- [13] Haroon Idrees, Imran Saleemi, Cody Seibert, and Mubarak Shah. Multi-source multi-scale counting in extremely dense crowd images. In *CVPR*, pages 2547–2554, 2013. 6
- [14] 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. 6
- [15] Wang Li, Li Yongbo, and Xue Xiangyang. Coda: Counting objects via scale-aware adversarial density adaption. In *ICME*, pages 193–198, 2019. 2
- [16] Chenxi Liu, Barret Zoph, Maxim Neumann, Jonathon Shlens, Wei Hua, Li-Jia Li, Li Fei-Fei, Alan Yuille, Jonathan Huang, and Kevin Murphy. Progressive neural architecture search. In *ECCV*, pages 19–34, 2018. 3
- [17] Yongtuo Liu, Qiang Wen, Haoxin Chen, Wenxi Liu, Jing Qin, Guoqiang Han, and Shengfeng He. Crowd counting via cross-stage refinement networks. *IEEE TIP*, 29:6800–6812, 2020. 1, 3
- [18] Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization. In *NeurIPS*, pages 7827–7838, 2018. 3, 5
- [19] Hyeonseob Nam, HyunJae Lee, Jongchan Park, Wonjun Yoon, and Donggeun Yoo. Reducing domain gap by reducing style bias. In *CVPR*, pages 8690–8699, 2021. 2
- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. 2014. 3, 6, 7
- [21] Vishwanath A Sindagi and Vishal M Patel. Inverse attention guided deep crowd counting network. In *AVSS*, pages 1–8, 2019. 5
- [22] Vishwanath A Sindagi, Rajeev Yasarla, Deepak Sam Babu, R Venkatesh Babu, and Vishal M Patel. Learning to count in the crowd from limited labeled data. In *ECCV*, pages 212–229, 2020. 2, 7
- [23] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *Int. Conf. Intell. Robots and Syst.*, pages 23–30, 2017. 2
- [24] Qi Wang, Junyu Gao, Wei Lin, and Yuan Yuan. Learning from synthetic data for crowd counting in the wild. In *CVPR*, pages 8198–8207, 2019. 2, 6, 7
- [25] Qi Wang, Junyu Gao, Wei Lin, and Yuan Yuan. Pixel-wise crowd understanding via synthetic data. *IJCV*, 129(1):225–245, 2021. 7
- [26] Qinwei Xu, Ruipeng Zhang, Ya Zhang, Yanfeng Wang, and Qi Tian. A fourier-based framework for domain generalization. In *CVPR*, pages 14383–14392, 2021. 2
- [27] Bin Yan, Houwen Peng, Kan Wu, Dong Wang, Jianlong Fu, and Huchuan Lu. Lighttrack: Finding lightweight neural networks for object tracking via one-shot architecture search. In *CVPR*, pages 15180–15189, 2021. 3
- [28] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In *CVPR*, pages 4085–4095, 2020. 3
- [29] Xiangyu Yue, Yang Zhang, Sicheng Zhao, Alberto Sangiovanni-Vincentelli, Kurt Keutzer, and Boqing Gong. Domain randomization and pyramid consistency: Simulation-to-real generalization without accessing target domain data. In *ICCV*, pages 2100–2110, 2019. 2
- [30] Cong Zhang, Kai Kang, Hongsheng Li, Xiaogang Wang, Rong Xie, and Xiaokang Yang. Data-driven crowd understanding: A baseline for a large-scale crowd dataset. *IEEE TMM*, 18(6):1048–1061, 2016. 6
- [31] Yingying Zhang, Desen Zhou, Siqin Chen, Shenghua Gao, and Yi Ma. Single-image crowd counting via multi-column convolutional neural network. In *CVPR*, pages 589–597, 2016. 6
- [32] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, pages 2223–2232, 2017. 7