

000  
001      **Unsupervised Deep Domain Adaptation on**  
002      **People Detection**  
003

004      Anonymous ECCV submission  
005

006      Paper ID \*\*\*  
007

009  
010      **Abstract.** This paper addresses the problem of unsupervised domain  
011 adaptation on the task of people detection in crowded scenes. That is,  
012 given a deep detection model well-trained on source domain, we adapt  
013 it into scene-specific detectors for any target domain on which no anno-  
014 tations are available. Firstly, we utilize iterative algorithm to iteratively  
015 auto-annotate target samples with high confidence on people instance as  
016 training set for scene-specific model on target domain. However, auto-  
017 annotated samples not only are lack of negative samples, but contains  
018 false positive samples. Therefore, on the one hand, we reuse negative  
019 samples from source domain to compensate for imbalance between the  
020 amount of positive samples and negative samples. On the other hand,  
021 we design an unsupervised regularizer based on deep network to mitigate  
022 influence from data error. Besides, we transform the last full connected  
023 layer into two sub-layers— element-wise multiply layer and sum layer, on  
024 which the unsupervised regularizer can be added on. In experiments on  
025 people detection, the proposed method boosts recall by nearly 30% while  
026 precision stays almost the same. Furthermore, we perform our method  
027 on standard domain adaptation benchmarks on both supervised and un-  
028 supervised settings and our results are state of the art.  
029

030  
031      **Keywords:** Unsupervised Domain Adaptation, Unsupervised Regular-  
032      izer, Deep Neural Network, People Detection  
033

034      **1 Introduction**  
035

036      Deep neural network has shown great power on traditional computer vision tasks,  
037 however, the labelled dataset should be large enough to train a deep model. In  
038 famous challenges such as PASCAL VOC and MS COCO, millions of labelled  
039 images are needed for training. This is also the case in surveillance applications.  
040 The annotation process for the task of people detection in crowded scenes is even  
041 more resource consuming, cause we need to label concrete locations of people  
042 instances. In modern society, there are over millions of cameras deployed for  
043 surveillance. However, these surveillance situations vary in lights, background,  
044 viewpoints, camera resolutions and so on. Directly utilizing models trained on  
old scenes will results in poor performance on the new situations due to data  
distribution changes. It is also unpractice to annotate people instances for every  
surveillance situation.

When there are few or even none of labelled data in target domain, domain adaptation helps to reduce the amount of labelled data needed. Most traditional works [1–5] either learn a shared representation between source and target domain, or project features into a common subspace. Recently, there are also works [6–8] proposed to learn a scene-specific detector by deep architectures. However, these approaches are heuristic either on constructing feature space or re-weighting samples. Our motivation of developing a domain adaptation architecture is to reduce heuristic methods required during adaptation process.

In this paper, we proposed a new approach of unsupervised deep domain adaptation on people detection. Using source model trained on source domain as initialization, we utilize iterative algorithm to iteratively auto-annotate target examples with high confidence as people instance on target domain for the first iteration. During each iteration, these auto-annotated data are regarded as training set to update target model, which, then, can be taken as the auto-annotation tool to auto-annotate target samples for the next iteration. However, these auto-annotated samples are defective, including lack of negative samples and existence of false positive samples, which will no doubt lead to exploration of predictions on non-people instances. Therefore, on the one hand, to compensate for the quantitative imbalance between positive and negative samples, we randomly sample negative instances from source domain and mix into training set. On the other hand, based on deep network, we design an unsupervised regularizer to mitigate influence from data error and avoid overfitting. To have better regularization effect during adaptation process, we transform the last full connected layer of deep model into two sub-layers, element-wise multiply layer and sum layer. Thus, the unsupervised regularizer can be added on element-wise multiply layer to adjust all weights in the deep network and gain better performance.

Also, we further evaluate our approach on standard domain adaptation benchmark Office Dataset. The results of our adaptation approach outperform previously published works on both supervised and unsupervised scenarios, which also demonstrate the feasibility of our adaptation approach on both detection and classification tasks.

The contributions of our work are three folds.

- We proposed a feasible scheme to learn scene-specific deep detectors for target domains by unsupervised methodology, which can be easily deployed to various surveillance situations without any additional annotations.
- For better performance of unsupervised regularizer, we transform the last full connected layer of deep network into two sub-layers, element-wise layer and sum layer. Thus, all weights contained in the deep network can be adjusted under the unsupervised regularizer. To our knowledge, this is the first attempt to transform full connected layers for the purpose of domain adaptation.
- Experiments on standard domain adaptation benchmarks for classification also demonstrate the applicability of our approach to other deep domain adaptation tasks.

The remainder of this paper is organized as follows. Section 2 reviews related works. Section 3 presents the details of our approach. Experimental results are shown in Section 4. Section 5 concludes the paper.

## 2 Relate Work

In many detection works, generic model trained by large amount of samples on source domain are directly utilized to detect on target domain. They assume that samples on target domain are subsets of source domain. However, when the distribution of data on target and source domain varies largely, the performance will drop significantly. Domain adaptation aims to reduce the amount of data needed for target domain.

Many domain adaptation works tried to learn a common representation space shared between source and target domain. Saenko et al. [1, 2] proposed a both linear-transform-based technique and kernel-transform-based technique to minimize domain changes. Gopalan et al. [3] projected features into Grassmann manifold instead of operating on features of raw data. Alternatively, Mesnil et al. [9] used transfer learning to obtain good representations. However, these methods are limited because scene-specific features are not learned to boost accuracy. The regularizer of our method are inspired these works.

Another group of works [4, 5, 10] on domain adaptation is to make the distribution of source and target domain more similar. Among these works, Maximum Mean Discrepancy (MMD) is used to as a metric to reselect samples from source domain in order to have similar distribution as target samples. In [11], MMD is incorporated as regularization to reduce distribution mismatch.

There are also works on deep adaptation to construct scene-specific detector. Wang et al.[6] explored context cues to compute confidence, and [7] learns distributions of target samples and proposed a cluster layer for scene-specific visual patterns. These works re-weighted auto-annotated samples for their final object function and additional context cues are needed for reliable performance. However, heuristic methods are required to re-weight samples. Alternately, Hattori el al. [8] learned scene-specific detector by generating a spatially-varying pedestrian appearance model. And Pishchulin et al. [12] used 3D shape models to generate training data. However, Synthesis for domain adaptation are also costly. Our approach minimize heuristic algorithms needed during adaptation process.

## 3 Our Approach

In this section, we introduce xxx. Iteration to auto-annotate. reasons for iteration: 1. unsupervised scenario. 2. avoid overfitting. problems in iteration, loss exceed gains. thus we added regularizer. regularizer is based on assumption xxx.

the adaptation architecture xxx. source stream xxx. target stream xxx. basic model xxx. transform xxx.

Denote xx as xxx..... we can define loss function as xxx. where xx is xxx.

Iteration algorithm.

separately training, in initialization stage, xxx. In adaptation stage xxx.

...

In this section, we introduce our unsupervised domain adaptation architecture on the task of people detection in crowded scenes. We denote training samples from source domain as  $\mathbf{X}^S = \{x_i^S\}_{i=1}^{N^S}$ . For training samples on source domain, we have corresponding annotations  $Y^S = \{y_i^S\}_{i=1}^{N^S}$  with  $y_i^S = (b_i^S, c_i^S)$ , where  $b_i^S = (x, y, w, h) \in R^4$  is the bounding box location and  $c_i^S \in \{0, 1\}$  is the label indicating whether  $x_i^S$  is a people instance. Also, we can denote the training samples on target domain as  $\mathbf{X}^{T,n} = \{x_j^{T,n}\}_{j=1}^{N^{T,n}}$  and corresponding annotations as  $Y^{T,n} = \{y_j^{T,n}\}_{j=1}^{N^{T,n}}$  with  $y_j^{T,n} = (b_j^{T,n}, c_j^{T,n})$  and  $c_j^{T,n} \equiv 1$ , where  $n \in \{1, \dots, N^I\}$  is the index of adaptation iteration process.  $N^I$  is the maximum number of adaptation iterations. Note that different from source domain, training set on target domain are auto-annotated samples with high confidence as people instance. Therefore, these auto-annotated data are all regarded as positive samples and equal amount of negative samples are randomly selected from target domain. During every adaptation iteration, the auto-annotation tool (also the target model) will be updated. Thus, the training samples for target domain at  $n^{th}$  adaptation iteration may differ from that at  $(n+1)^{th}$  adaptation iteration.

The adaptation architecture of our approach consists of two streams – source stream and target stream, as shown in Fig 1. Source steam takes samples from source domain as input, while target stream are trained by auto-annotated positive samples from target domain and negative samples from source domain. These two streams can utilize any deep detection network as their basic model, as well as their detection loss function as supervised loss functions of two streams. In our experiment, we use the detection network mentioned in Section 3.3 as the basic model. For the purpose of unsupervised regularizer, the last full connected layers on both source and target model are transformed into two sub-layers – element-wise multiply layer and sum layer, as mentioned in Section 3.2.

These two streams are separately trained at different stage of adaptation process. At initialization stage, source model of source stream are trained under supervised loss function with abundant labelled data,  $\mathbf{X}^S$ , from source domain. After its convergence, the weights of source model  $\theta^S$  are taken to initialize target stream. At adaptation stage, target model is trained by iteration algorithm stated in Section 3.1 under both supervised loss function and unsupervised regularizer. Note that we do not jointly train two streams at adaptation stage and the weights of source model stays static which is served as a distribution reference for unsupervised regularizer at adaptation stage.

The loss function of adaptation architecture is composed of a supervised loss and an unsupervised regularizer. To learn a scene-specific detector for the target domain, we directly use the original detection loss function from basic detection network as the supervised loss in our architecture. As the auto-annotated samples contain data error, the performance decline caused by false positive samples will exceed the performance boost resulting from true positive samples if the target

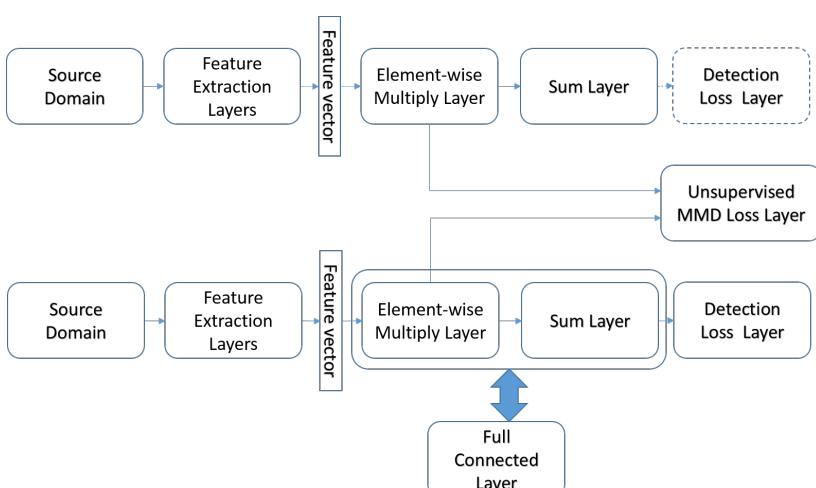
model is trained only with supervised loss function. Therefore, an unsupervised regularizer is required to mitigate the influence from data error on target model. Based on the assumption that source domain and target domain should share the same feature space after feature extraction layers, we encode an unsupervised regularizer to make a constraint that the distribution of data representation on the element-wise multiply layer should be similar between source stream and target stream. We formulate the combination of supervised loss and unsupervised regularizer as follows:

$$L(\theta^{T,n} | \mathbf{X}^{T,n}, \mathbf{Y}^{T,n}, \mathbf{X}^S, \theta^S) = L_S + \alpha * L_U \quad (1)$$

$$\begin{aligned} L_S &= \sum_{j=1}^{N^{T,n}} (R(\theta^{T,n} | x_j^{T,n}, b_j^{T,n}) + C(\theta^{T,n} | x_j^{T,n}, c_j^{T,n})) \\ &\quad + \sum_{i=f(j)}^{N^{T,n}} (R(\theta^{T,n} | x_i^S, b_i^S) + C(\theta^{T,n} | x_i^S, c_i^S)) \end{aligned} \quad (2)$$

$$L_U = L_{EWM}(\theta^T | \mathbf{X}^T, \mathbf{X}^S, \theta^S) \quad (3)$$

where  $L_S$  is supervised loss to learn the scene-specific detector and  $L_U$  is the unsupervised regularizer part.  $R(\cdot)$  is a regression loss for bounding box location, like norm-1 loss, and  $C(\cdot)$  is a classification loss for bounding box confidence, such as cross-entropy loss.  $f(\cdot)$  randomly returns an index under the condition that  $x_{f(\cdot)}^S$  must be a negative sample. And  $L_{EWM}(\cdot)$ , to be introduced in Section 3.2, is the MMD-based loss added on the element-wise multiply layer for unsupervised regularization.  $\alpha$  is the coefficient balancing the effect of supervised and unsupervised loss.



**Fig. 1.** xxxxxxxxxxxxxxxxxxxxxxxxx

### 225 3.1 Iterative Algorithm

226 In this section we introduce the iteration algorithm as training method at adap-  
 227 tation stage. The auto-annotating tool takes images on target domain with high  
 228 confidence as training data for next iteration. To generate the auto-annotated  
 229 data for the first iteration, we utilize the generic model well-trained on source  
 230 domain, which results in training set on the target domain. As mentioned in  
 231 Sec 3.3, the training set on target domain auto-annotated by the generic model  
 232 in our experiment have low recall and high precision. At subsequent iterations,  
 233 training set on target domain are auto-annotated by upgraded model resulting  
 234 from training of last iteration. Among every iteration, these auto-annotated data  
 235 are used as training data to upgrade the deep network.

236 There are two reasons to employ iteration algorithm. Firstly, auto-annotated  
 237 data on target domain change for every adaptation iteration and new samples  
 238 will be auto-annotated as positive instances. Compared to methods without it-  
 239 eration algorithm, it helps to avoid overfitting caused by lack of data. Besides,  
 240 unsupervised regularizer performs better with more training data as it's a dis-  
 241 tribution based regularizer.

242 As the training set on target domain auto-annotated for the training of first  
 243 iteration have low recall rate and people and non-people regions mixed in nega-  
 244 tive instances, we ignore back propagation of bounding boxes with low confidence  
 245 during training. That is, we encourage the network to have more confidence on  
 246 positive instances and stay conservative toward negative instances. This policy  
 247 will no doubt resulting predictions of many non-people instance and impeding  
 248 the further training when many non-people instances are regarded as people in-  
 249 stances by the target model. To compensate for lack of true negative instances,  
 250 we add annotated samples from source domain into the training data, which  
 251 are human annotated and can thus provide true negative samples. In our ex-  
 252 periment, when training on additional samples from source domain, we only do  
 253 back propagation of bounding boxes with low confidence. At the same time,  
 254 the unsupervised loss will also regularize the network. The complete adaptation  
 255 process is illustrated in Algorithm 1. After a predetermined iteration limit  $N^I$   
 256 is reached, we obtain our final detection model on the target domain.

### 258 3.2 Unsupervised weights regularizer on Element-wise Multiply 259 Layer

261 **Element-wise Multiply Layer** In deep neural network, the last feature vector  
 262 layer are taken as an important data representation of input images. However,  
 263 in this paper, we take one step further to focus on the last full connected layer  
 264 which serves as an decoder to decode rich information contained in the last  
 265 feature vector into final outputs. As source model are trained with abundant  
 266 labelled data on source domain, the parameters of the last full connected layer  
 267 are also well converged. We assume that a regularizer on the last full connected  
 268 layer will achieve better results compared with the last feature vector layer.  
 269 Firstly, denote the last feature vector, parameters of the last full connected layer

---

**Algorithm 1** Deep domain adaptation algorithm (to be completed)

---

```

270 1: procedure DEEP DOMAIN ADAPTATION
271 2: Train source model on source stream with abundant annotated data
272 3: Use well-trained source model on source stream to initialize model on target stream
273 as  $M_0$ 
274 4: for  $i = 0:N^I$  do
275 5:    $M_i$  generate "fake ground truth"  $G_i$  of target domain
276 6:   balbal
277 7:   balbabla
278 8:    $G_i = G_i +$  random samples from source domain
279 9:   Take  $G_i$  as training data to upgrade  $M_i$  into  $M_{i+1}$ 
280 10: end for
281 11:  $M_{N^I}$ : final model.
282 12: end procedure

```

---

and final outputs as  $\mathbf{F}_{(N^B \times N^D)}$ ,  $\mathbf{C}_{(N^D \times N^O)}$  and  $\mathbf{P}_{(N^B \times N^O)}$ , respectively. The operation of full connected layer can be thus formulated as matrix multiply:

$$\mathbf{P} = \mathbf{F} * \mathbf{C} \quad (4)$$

$$P_{b,o} = \sum_d F_{b,d} * K_{d,o} \quad (5)$$

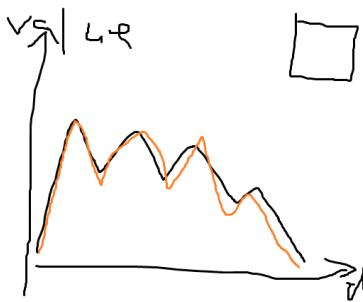
Inspired by this form, we separate the above formula into two sub-operations – element-wise multiply and sum, which can be formulated as:

$$M_{b,o,d} = F_{b,d} * C_{d,o} \quad (6)$$

$$P_{b,o} = \sum_d M_{b,o,d} \quad (7)$$

where  $\mathbf{M}_{(N^B \times N^O \times N^D)} = [\mathbf{m}_{b,o}]$  is the parameter tensor of element-wise multiply operations.  $\mathbf{m}_{b,o}$  is a vector with  $N^D$  dimensions, which will be the basis of unsupervised regularizer. Finally, we can equivalent-transform the last full connected layer between the last feature vector layer and final outputs layer into element-wise multiply layer and sum layer. The transformed element-wise multiply layer is thus the last layer with parameters before output layers. Fig x (!) illustrates the transform.

**Unsupervised weights regularizer on Element-wise Multiply Layer** In works [decaf][], the last feature vector layer are regarded as the final representation of images. (how to introduce beyond sharing weights) In domain adaptation tasks, when generic deep model are trained with abundant data from source domain, the last element-wise multiply layer mentioned in Sec 3.2 also includes rich information leading to the final output. For the nodes in the output layer, inputs from element-wise multiply layers that will contribute to its value are not randomly decided. In this paper, we assume that the distribution of  $\mathbf{m}_{b,o}$  of



**Fig. 2.** One kernel at  $x_s$  (*dotted kernel*) or two kernels at  $x_i$  and  $x_j$  (*left and right*) lead to the same summed estimate at  $x_s$ .

the last element-wise multiply layer on both source and target domain should be similar. In the last element-wise multiply layer, every dimension of  $\mathbf{m}_{b,o}$  may contribute to the final output. However, as the old model on source domain are trained with abundant images, the distribution of  $\mathbf{m}_{b,o}^T$  should be consistent compared with  $\mathbf{m}_{b,o}^S$  of source domain when adapting deep network on target domain. We utilize MMD (maximum mean discrepancy) to encode the similarity of element-wise multiply layer of source domain and target domain:

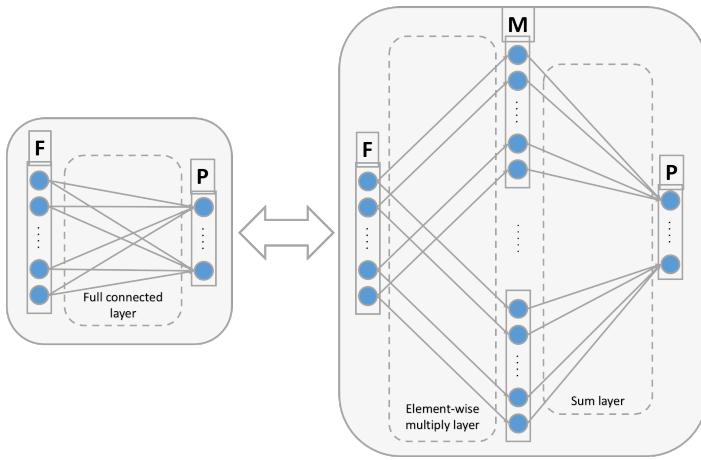
$$L_{MMD}(\theta^T | \mathbf{X}^S, \mathbf{X}^T, \theta^S) = \frac{1}{N^O} \sum_{o=1}^{N^O} \left\| \frac{1}{N^B} \sum_{b=1}^{N^B} \mathbf{m}_{b,o}^T - \frac{1}{N^B} \sum_{b=1}^{N^B} \mathbf{m}_{b,o}^S \right\|^2 \quad (8)$$

which can also be interpreted as the Euclidean distance between the center of  $\mathbf{m}_{b,o}^T$  and  $\mathbf{m}_{b,o}^S$  across all output nodes. It's unpractical to get the distribution of the whole training set, while too few images cannot obtain a stable center for regularization. In our experiments, the  $L_{MMD}(\cdot)$  loss is calculated for every batch. An example comparison of centers of  $\mathbf{m}_{b_i,o}^S$  and  $\mathbf{m}_{b_j,o}^S$  are shown in Fig x(1).

### 3.3 Detection Network

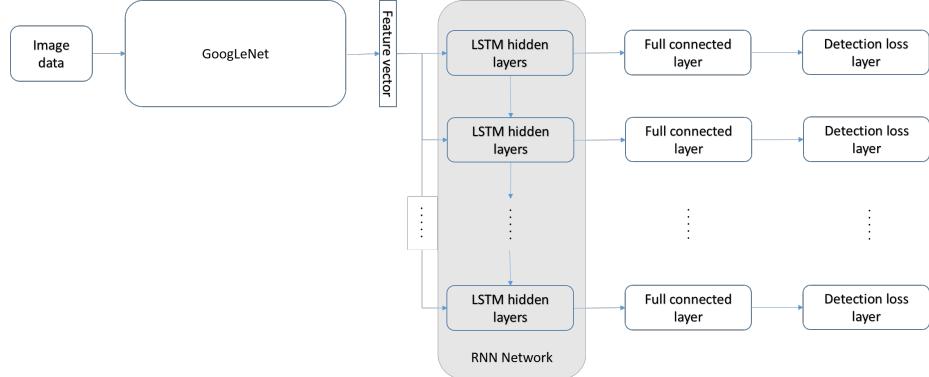
The generic model <sup>1</sup> used in our adaptation architecture for source and target stream is an end to end detection network without any precomputed region proposals needed. It consists of a GoogLeNet [13] for feature extraction and a RNN-based decoder for output of bounding box and corresponding confidence, as shown in Fig x(!). Firstly, the GoogLeNet model encode the image into a feature map (15x20x1024) of which each 1024 dimension vector are data representation of its receptive field corresponding to a subregion of the image. Then the RNN-based layers with batch size 300 will decode the data representation

<sup>1</sup> Proposed by Russel et al.



**Fig. 3.** xxxxxxxxxxxxxxxxx

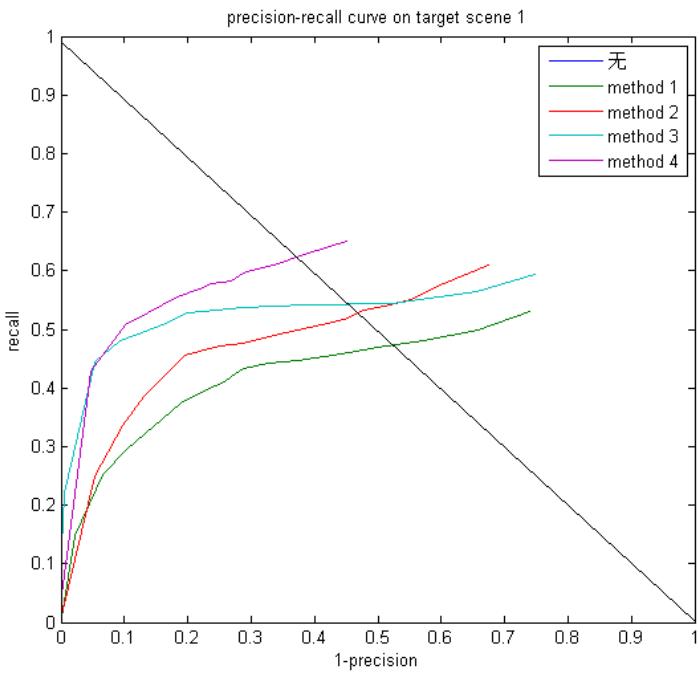
and sequentially predict 5 possible bounding boxes by the order of its corresponding confidence. Finally, all outputs are summarized to give final detection results. When trained with abundant samples from source domain, the obtained model has a high precision on target domain, however, its recall is low. Also, different from other detection network [14, 15], which need precomputed proposals for classification and fine regression, this generic model directly predict bounding boxes with high confidence. Thus, negative instances may contain people and non-people predictions, and cannot be employed in adaptation training.



**Fig. 4.** xxxxxxxxxxxxxxxxx

## 4 Experiment Results

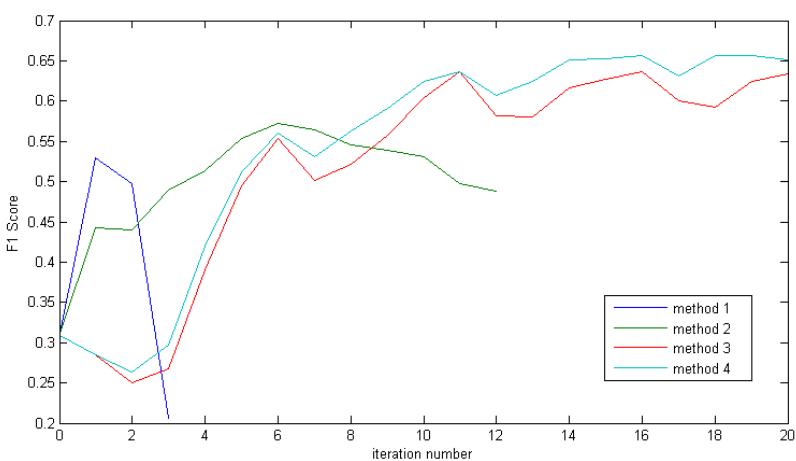
In this section, we introduce experiment results on both surveillance applications and standard domain adaptation dataset. Our motivation for unsupervised domain adaptation method is for easier deployment of We firstly evaluate our approach on video surveillance. Then we employ our approach to standard domain adaptation benchmarks on both supervised and unsupervised settings to demonstrate the effectiveness of our method.



**Fig. 5.** xxxx

### 4.1 Domain Adaptation on Crowd Dataset

**Dataset and evaluation metrics** To show the effectiveness of our domain adaptation approach on people detection, we collected a dataset consisting of 3 target scenes for target domain. These three scenes contain 1308, 1213 and 331 un-annotated images with 0000, 0000, 0000 people instances respectively. For each scene, 100 images are annotated for evaluation. Instead of labelling the whole body of a person, we labels the head of a person as bounding box during training. The motivation for labelling only people heads comes from detection of

**Fig. 6.** xxxx

indoor people or in crowded scenes, where the body of a person may be invisible. The dataset for source domain are Brainwash Dataset released on <http://d2.mpi-inf.mpg.de/datasets>. Brainwash Dataset consists of over 11917 images from 3 crowded scenes. Examples of images from source and target domain are shown in Fig x(!).

Our evaluation metrics for detection uses the protocol defined in PASCAL VOC [16]. To judge a predicted bounding box whether correctly matches a ground truth bounding box, their intersection over their union must exceed 50%. And Multiple detections of the same ground truth bounding box are regarded as one correct prediction. We plot the precision-recall curve in Fig x(!). Also, the F1 score  $F1 = 2 * precision * recall / (precision + recall)$  during the adaptation process are also shown in Fig x(!).

**Experimental settings** We use deep learning framework Caffe [17] as the adaptation architecture of our approach. During the adaptation, we set learning rate as 0.01 and momentum as 0.5. At initialization stage, GoogLeNet weights are firstly used to initialize source model of source stream, while parameters in RNN layers are randomly initialized from a uniform distribution. For each iteration, 100 auto-annotated images from target domain and 1000 annotated images from source domain are alternatively used for training. The outputs of our detection network contains bounding box locations and corresponding confidence, thus there are two full connected layers between the last feature vector layer and the final outputs. Experiments of unsupervised regularizer on element-wise multiply layer for bounding box regression have no additional improvement on the performance when regularizer on element-wise multiply layer for box confi-

dence classification are added already. Our approach on domain adaptation are executed separately in the 3 target scenes.

**Comparison with baseline methods** To demonstrate the effectiveness of our approach, 4 methods are compared with method 4 as our final approach:

1. Only auto-labeled samples on target domain are used for training, and without any unsupervised regularizer.
2. Only auto-labeled samples on target domain are used for training, with M-MD regularizer on last element-wise multiply layer as unsupervised weights regularizer.
3. Both auto-labeled images from target domain and labeled images from source domain are alternately sampled for training, with MMD regularizer on last feature vector as unsupervised weights regularizer.
4. Both auto-labeled images from target domain and labeled images from source domain are alternately sampled for training, with MMD regularizer on last element-wise multiply layer as unsupervised weights regularizer.

Fig x(!) plots the precision-recall curve of the above comparison methods in target scene 1). Also, the F1 score changes of every iteration during adaptation process are also depicted in Fig x(!). Table x(!) gives concrete precision and recall value of the 4 comparison methods on three target scenes when the F1 scores are at their highest. Examples of adaptation results are shown in Fig x(!).

	Scene 1			Scene 2			Scene 3		
	1-Pr	Re	F1	1-Pr	Re	F1	1-Pr	Re	F1
method 0	0.101	0.187	0.309	0.015	0.683	0.807	0.035	0.412	0.577
method 1	0.245	0.408	0.530	0.632	0.905	0.524	0.176	0.778	0.800
method 2	0.284	0.476	0.572	0.012	0.837	<b>0.906</b>	0.078	0.653	0.764
method 3	0.109	0.496	0.637	0.002	0.721	0.838	0.044	0.611	0.746
method 4	0.140	0.530	<b>0.656</b>	0.006	0.811	0.893	0.097	0.778	<b>0.836</b>

**Performance evaluation** From the Table x(!) and Fig x(1), we have the following observations:

- The recall values of method 1,2,3,4, which all utilized iteration algorithm to upgrade the target model, are larger than that of method 0, which are source model trained on source domain. This implies the effectiveness of our iteration algorithm in auto-annotation and iterative training.
- Compared with method 1, method 2 has higher F1 score. Their difference on whether a MMD regularizer are added into loss function demonstrates that our unsupervised regularizer can suppress data error and thus boost the final recall.

- 540 – Method 4 has both higher precision and higher recall than method 2, which  
 541 demonstrates the effectiveness of additional samples from source domain  
 542 during adaptation process.  
 543 – Compared with method 3, the recall of method 4 are further boosted. This  
 544 results from the transformed element-wise multiply layer which provided  
 545 better regularization effect on the target model.

## 4.2 Domain Adaptation on Standard Classification Benchmark

**Office dataset** The Office dataset [1] comprises 31 categories of objects from 3 domains (Amazon, DSLR, Webcam). Example images are depicted in Fig. x(!). As Amazon domain contains 2817 labelled images, which is the largest, we take it as source domain and Webcam domain as target domain. We follow the standard protocol for both supervised and unsupervised settings. Specifically, for supervised domain adaptation, we use 20 randomly sampled images with labels for each category as training data for Amazon domain. When evaluate on unsupervised domain adaptation, 3 labelled images from target domain are additional selected for each class. For both settings, the rest of images on target domain are used for evaluation.



559  
 560  
 561 **Fig. 6.** Some examples from three domains in the Office dataset. |  
 562  
 563  
 564  
 565  
 566  
 567

571 **Fig. 7.** One kernel at  $x_s$  (dotted kernel) or two kernels at  $x_i$  and  $x_j$  (left and right)  
 572  
 573  
 574

**575 Experimental settings and network design** On supervised setting, we  
 576 reused the architecture in people detection. We utilize AlexNet [18] as the generic  
 577 model of both streams. Firstly, we train the source model on source stream  
 578 with provided training data from Amazon domain. Then iteration algorithm  
 579 mentioned in Sec 3.1 are utilized for adaptation. The auto-labelling tool takes  
 580 images on target domain with confidence exceeding 0.9 as training data for next  
 581 iteration. The difference is besides auto-labelled images on target domain, 3 human  
 582 labelled images for each class are also included as training data for target  
 583 model. For each iteration, 100 images are randomly sampled from training data.  
 584 The unsupervised MMD regularizer added on the element-wise multiply layer

transformed from the last full connect layer of target model set the coefficient value  $\alpha$  as 10 on Eq 3, the same as that in people detection task.

We use the same experimental setting for our unsupervised adaptation, except that at adaptation stage, no human labelled images can be added into the training set.

**Performance evaluation** In Table x(!), we compare our approach with other six recently published works in both supervised and unsupervised settings. The outstanding performance on both settings confirms the effectiveness of our iteration algorithm and MMD regularizer on the element-wise multiply layer transformed from the last full connect layer.

	$A \rightarrow W$	
	Supervised	Unsupervised
GFK(PLS,PCA)[19]	46.4	15.0
SA [20]	45.0	15.3
DA-NBNN [21]	52.8	23.3
DLID [22]	51.9	26.1
DeCAF <sub>6</sub> S [23]	80.7	52.2
DaNN [11]	53.6	35.0
Ours	<b>84.3</b>	<b>66.3</b>
Ours	<b>85.4</b>	<b>69.3</b>

## 5 Conclusions

The paper ends with a conclusion.

## References

1. Saenko, K., Kulis, B., Fritz, M., Darrell, T.: Adapting visual category models to new domains. In: Computer Vision–ECCV 2010. Springer (2010) 213–226
2. Kulis, B., Saenko, K., Darrell, T.: What you saw is not what you get: Domain adaptation using asymmetric kernel transforms. In: Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on, IEEE (2011) 1785–1792
3. Gopalan, R., Li, R., Chellappa, R.: Domain adaptation for object recognition: An unsupervised approach. In: Computer Vision (ICCV), 2011 IEEE International Conference on, IEEE (2011) 999–1006
4. Huang, J., Gretton, A., Borgwardt, K.M., Schölkopf, B., Smola, A.J.: Correcting sample selection bias by unlabeled data. In: Advances in neural information processing systems. (2006) 601–608
5. Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K., Schölkopf, B.: Covariate shift by kernel mean matching. Dataset shift in machine learning **3**(4) (2009) 5



Fig. 8. One kernel at  $x_s$  (dotted kernel) or two kernels at  $x_i$  and  $x_j$  (left and right)

- 675 6. Wang, X., Wang, M., Li, W.: Scene-specific pedestrian detection for static video  
676 surveillance. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*  
677 **36**(2) (2014) 361–374
- 678 7. Zeng, X., Ouyang, W., Wang, M., Wang, X.: Deep learning of scene-specific clas-  
679 sifier for pedestrian detection. In: *Computer Vision–ECCV 2014*. Springer (2014)  
680 472–487
- 681 8. Hattori, H., Naresh Boddeti, V., Kitani, K.M., Kanade, T.: Learning scene-specific  
682 pedestrian detectors without real data. In: *Proceedings of the IEEE Conference*  
683 *on Computer Vision and Pattern Recognition*. (2015) 3819–3827
- 684 9. Mesnil, G., Dauphin, Y., Glorot, X., Rifai, S., Bengio, Y., Goodfellow, I.J., Lavoie,  
685 E., Muller, X., Desjardins, G., Warde-Farley, D., et al.: Unsupervised and transfer  
686 learning challenge: a deep learning approach. *ICML Unsupervised and Transfer*  
687 *Learning* **27** (2012) 97–110
- 688 10. Gong, B., Grauman, K., Sha, F.: Connecting the dots with landmarks: Discrimi-  
689 natively learning domain-invariant features for unsupervised domain adaptation.  
690 In: *Proceedings of The 30th International Conference on Machine Learning*. (2013)  
222–230
- 691 11. Ghifary, M., Kleijn, W.B., Zhang, M.: Domain adaptive neural networks for object  
692 recognition. In: *PRICAI 2014: Trends in Artificial Intelligence*. Springer (2014)  
898–904
- 693 12. Pishchulin, L., Jain, A., Wojek, C., Andriluka, M., Thormählen, T., Schiele, B.:  
694 Learning people detection models from few training samples. In: *Computer Vision*  
695 and *Pattern Recognition (CVPR), 2011 IEEE Conference on*, IEEE (2011) 1473–  
696 1480
- 697 13. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D.,  
698 Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: *Proceedings*  
699 *of the IEEE Conference on Computer Vision and Pattern Recognition*. (2015) 1–9
- 700 14. Girshick, R.: Fast r-cnn. In: *Proceedings of the IEEE International Conference on*  
701 *Computer Vision*. (2015) 1440–1448
- 702 15. Vu, T.H., Osokin, A., Laptev, I.: Context-aware cnns for person head detection.  
703 In: *Proceedings of the IEEE International Conference on Computer Vision*. (2015)  
2893–2901
- 704 16. Everingham, M., Eslami, S.A., Van Gool, L., Williams, C.K., Winn, J., Zisser-  
705 man, A.: The pascal visual object classes challenge: A retrospective. *International*  
706 *Journal of Computer Vision* **111**(1) (2015) 98–136
- 707 17. Jia, Y., Shelhamer, E., Donahue, J., Karayev, S., Long, J., Girshick, R., Guadarrama,  
708 S., Darrell, T.: Caffe: Convolutional architecture for fast feature embedding.  
709 arXiv preprint arXiv:1408.5093 (2014)
- 710 18. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep  
711 convolutional neural networks. In: *Advances in neural information processing systems*.  
712 (2012) 1097–1105
- 713 19. Gong, B., Shi, Y., Sha, F., Grauman, K.: Geodesic flow kernel for unsupervised  
714 domain adaptation. In: *Computer Vision and Pattern Recognition (CVPR), 2012*  
715 *IEEE Conference on*, IEEE (2012) 2066–2073
- 716 20. Fernando, B., Habrard, A., Sebban, M., Tuytelaars, T.: Unsupervised visual do-  
717 main adaptation using subspace alignment. In: *Proceedings of the IEEE Interna-*  
718 *tional Conference on Computer Vision*. (2013) 2960–2967
- 719 21. Tommasi, T., Caputo, B.: Frustratingly easy nbnn domain adaptation. In: *Proce-*  
720 *dings of the IEEE International Conference on Computer Vision*. (2013) 897–904

- 720 22. Chopra, S., Balakrishnan, S., Gopalan, R.: Dlid: Deep learning for domain adap-  
721 tation by interpolating between domains. In: ICML workshop on challenges in  
722 representation learning. Volume 2. (2013)
- 723 23. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.:  
724 Decaf: A deep convolutional activation feature for generic visual recognition. arXiv  
725 preprint arXiv:1310.1531 (2013)