

Learning a Domain-Invariant Embedding for Unsupervised Person Re-identification

Nan Pu
LIACS Media Lab
Leiden University
Leiden, The Netherlands
n.pu@liacs.leidenuniv.nl

Georgiou T.K
LIACS Media Lab
Leiden University
Leiden, The Netherlands
t.k.georgiou@liacs.leidenuniv.nl

Erwin M. Bakker
LIACS Media Lab
Leiden University
Leiden, The Netherlands
erwin@liacs.nl

Michael S. Lew
LIACS Media Lab
Leiden University
Leiden, The Netherlands
m.s.k.lew@liacs.leidenuniv.nl

Abstract—Person re-identification (Re-ID) aims at matching images of the same person where images are captured by non-overlapping camera views distributed at different locations. To solve this problem, most recent works require a large pre-labeled dataset for training a deep model. These methods are not always suitable for real-world applications, because the latter often lack labeled data. In order to tackle this drawback, we proposed a novel Domain-Invariant Embedding Network (DIEN) to learn a domain-invariant embedding (DIE) feature by introducing a multi-loss joint learning with Recurrent Top-Down Attention (RTDA) mechanism. Due to the improvement in traditional triplet loss, our proposed model can benefit from both source-domain (labeled) data and target-domain (unlabeled) data. Furthermore, the resulting DIE feature not only has improved class discrimination but also robustness to domain shift. We compared our method with recent competitive algorithms and also evaluated the effectiveness of the proposed modules.

Index Terms—Person Re-identification, Unsupervised Domain Adaptive, Attention Mechanism

I. INTRODUCTION

In recent years, Person re-identification (Re-ID) in large-scale surveillance systems has become one of the most challenging and hottest topics of computer vision. The Re-ID technology helps us to match pedestrian images which include the same person and are captured by different cameras at different locations or the same cameras at different time.

In order to overcome various changes in appearance and environment, current deep learning based Person Re-ID models are focused in learning robust features automatically (e.g. end-to-end learning) instead of handcrafted features. Most existing Re-ID works focus on supervised methods. They utilize deep CNNs [1–4] to boost the performance. Nevertheless, these methods achieve significant performance improvements only when a large amount of labeled training data is available. In real Re-ID scenarios, by using mature pedestrian detection technology, we can conveniently obtain very large Re-ID datasets without labels [5]. Since labeling data is a very time consuming procedure, there are only a small number of datasets for these methods to be trained on. So, if we can transfer the Re-ID capability of a deep neural network, trained on a fully labeled dataset, to perform Re-ID on another unlabeled dataset, we may make accurate Re-ID more tractable. Usually, related works treat two different

datasets as a source domain (fully labeled) and target domain (without label). It is well known that Re-ID models trained on one domain often fail to generalize well to another [6]. Some researchers handle this problem by utilizing Unsupervised Domain Adaption (UDA) method [6–11]. And other works treat it as a transfer learning problem [12]. Both approaches need to make use of the unlabeled data to alleviate this drawback. In general, we can regard this problem as a domain shift or dataset shift [13] and thus apply a domain adaptive method to solve it. However, in unsupervised domain adaptive person Re-ID the two datasets do not share class labels (person identity), which is different from in traditional domain adaptive method. The challenge lies on how to obtain semantically meaningful domain-invariant features with good robustness for each identity. That is the main goal of our work presented in this paper.

To address above mentioned problems, we proposed a Domain-Invariant Embedding Network (DIEN), taking advantage of both source-domain (labeled) data and target-domain (unlabeled) data by using a novel proposed centering constrained cross-domain triplet loss (CCCDTL) function, to learn a domain-invariant embedding (DIE) feature for cross-domain Person Re-ID. Due to the supervision of the source domain and the auxiliary information of the target domain, the DIE feature is not only very discriminative, but is also robust under domain shift.

In order to further improve the discriminative power of DIE feature and the supervised information propagation, we introduce a new Recurrent Top-Down Attention (RTDA) module to recurrently find the region of interest on feature maps and re-weight each channel of the feature maps to enable knowledge distillation. This is achieved by multi-loss joint learning and iteratively updating the parameters of the attention module. After finishing DIE feature learning, our model can perform cross-domain Re-ID by directly retrieving DIE features of the query image and the gallery images.

The main novel contributions of our paper can be summarized as follows:

- We propose a novel centering constrained cross-domain triplet loss (CCCDTL) function to achieve cross-domain learning. By using this loss function, our model can make full use of labeled and unlabeled data simultaneously.

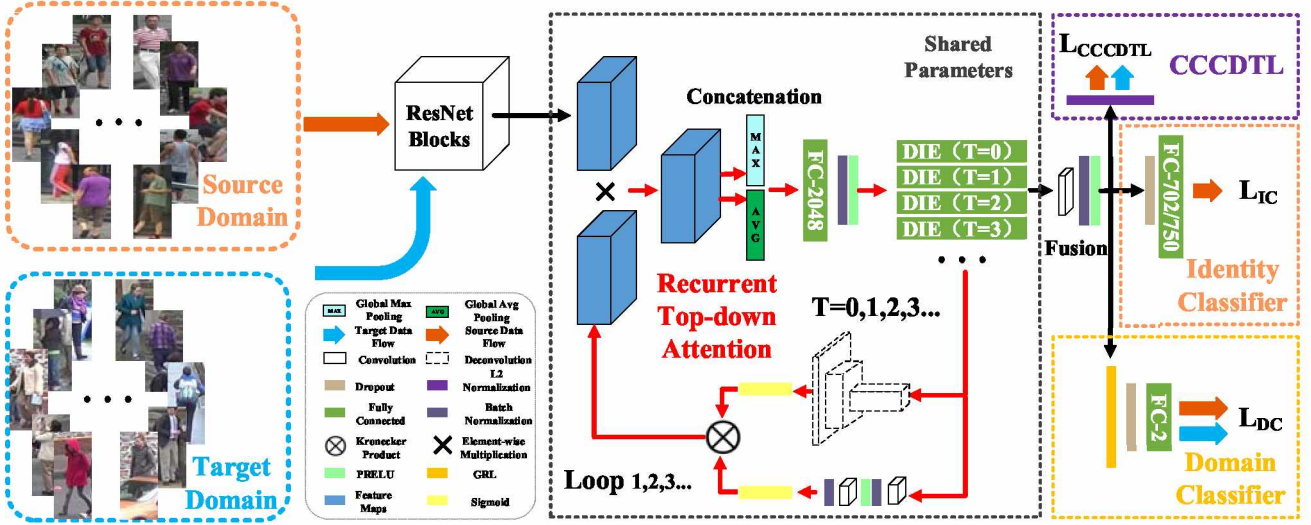


Fig. 1. Illustration of the deep domain-invariant embedding neural network (DIEN). The proposed model consists of a backbone network, an identity classifier (IC), a centering constrained cross-domain triplet loss (CCCDTL) function, a domain classifier (DC) and a Recurrent Top-Down Attention (RTDA) module.

- Our proposed DIE Network (DIEN) is a new end-to-end deep domain adaptive model. It is capable of learning domain-invariant embedding (DIE) features and recurrently refine learned feature by the Recurrent Top-Down Attention (RTDA) module proposed in this paper.

II. RELATED WORK

A. Person Re-Identification (Re-ID)

Supervised Learning for Re-ID: Most existing Re-ID models [1–3] are trained using supervised learning strategies. For example, in order to handle body parts misalignment, Suh et al. [1] propose a two-stream network to learn a part-aligned representation for person Re-ID by using a bilinear-pooling layer. Further, He et al. [2] present a Deep Spatial feature Reconstruction (DSR) method to address the partial person Re-ID problem. Recently, Conditional Random Fields (CRFs) are exploited to mine second-order relationships of mini-batch training data in [3], which dramatically improves the performance of deep neural networks for Re-ID. Although those methods achieve a significant increasing performance on recent datasets, namely Market-1501 [14] and DukeMTMC-ReID [15], these methods may not be practical since collecting a large amount of annotated training data depends on lots of manpower and time.

Unsupervised Learning for Re-ID: To alleviate the above limitation, researchers also focus on person Re-ID using unlabeled training data [9, 16]. As an example, Li et al. [16] take full advantage of the information of cameras in the target domain, treating multiple one-person images from different cameras as a tracklet. Another typical work introduces a progressive unsupervised learning (PUL) method [9], which utilized a clustering method to select representative samples to modify the pre-trained model. PUL aimed at transferring pre-trained deep representations to an unseen domain by a

Self-paced Learning. Nevertheless, due to the lack of label information for images across different cameras, unsupervised learning based methods typically can not perform as well as the supervised methods do.

B. Unsupervised Domain Adaptation for Re-ID

Unsupervised domain adaptation (UDA) has been studied widely in various computer vision tasks [10, 13, 17] and recently faces new challenges in person re-ID [8].

From a UDA perspective, most related works are concentrated around Domain-Invariant Feature Learning [10]. Some recent works leverage an auto-encoder to achieve knowledge distillation [11] so as to learn a domain-invariant representation with significant generalization. In order to increase reasonable cues for person Re-ID and decrease the influence of camera variance, Zhong et al. utilize CamStyle [7] to generate camera-style images [8] as extra training data. In [18], Wang et al. employed both the attribute and identity labels to encode an embedding feature to promote unsupervised cross-dataset or cross-domain Re-ID. Due to the success of generation models, now many cross-domain tasks are dominated by GAN-based methods such as Similarity Preserving Generative Adversarial Network (SPGAN) [6] and Person Transfer Generative Adversarial Network (PTGAN) [12]. Both of them showed that using data augmentation methods strengthens the Re-ID ability of a deep neural network on the target domain thus improving the performance and closing the domain gap.

In this paper, we follow the general setting of unsupervised domain adaptation as used in [8]. Specifically, we provide labeled source training images and unlabeled target training images as training data and evaluate the performance of the proposed model on the target testing database.

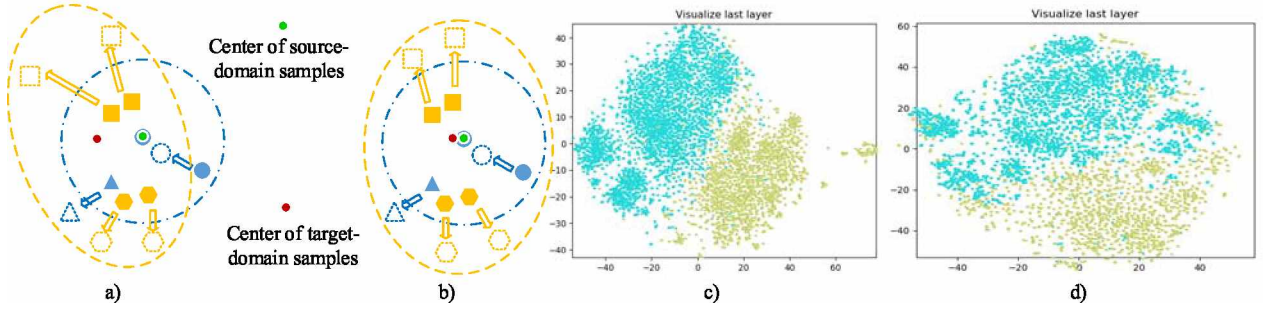


Fig. 2. Graph a) and b) represents different effect of Triple loss in [10] and our CCCDTL. The t-SNE of the pre-train features and learned DIE features are drawn in scatter c) and d) respectively (blue and green point denotes target and source data respectively).

C. Attention Mechanism

It is well known that attention plays an important role in human perception [19]. Recently, there have been several attempts [4, 20] to incorporate attention processing, to improve the performance of CNNs in Person Re-ID tasks. For example, Li et al. [20] designed a new two-stream model to learn both local and global features by hard and soft attention interactive learning. By refining the feature maps, their network not only performs well but is also robust to noisy inputs.

Unlike the attention-based methods in [4, 20], our proposed Recurrent Top-Down Attention (RTDA) leads deep neural network to recurrently update parameters and take the high-level feedback signal into feature extraction instead of extracting feature vectors based on one-pass of the data through the network forward.

III. PROPOSED METHOD

A. Network Architecture Overview

In our deep domain-invariant embedding neural network, we deploy five blocks of ResNet-50 [21] as a primary feature extractor and follow the training strategy in [8] which fine-tunes on the ImageNet pre-trained model. We construct our backbone network by replacing Global Average Pooling (GAP) layer and the last 1,000-dim fully connected (FC) layer with two pooling layers and a 2,048-dim FC layer followed by batch normalization (BN) [22] and PReLU [23], as shown in Fig.1.

Inspired by CBAM [24], we use both average-pooled and max-pooled features simultaneously to keep the distinctive object clues gathered by max-pooling. Specifically, we concatenate the two outputs of the global max pooling and global average pooling and feed them to the next FC layer. The output of this FC layer is a 2,048-dim feature vector, which we call the “domain-invariant embedding” (DIE).

For the purposes of strengthening the information flow and distilling the DIE feature, Recurrent Top-Down Attention (RTDA) is exploited to recurrently re-weight the channel and spatial position of feature maps simultaneously. The RTDA module is implemented by multiple deconvolution and convolution layers whose details will be described in the Section III-D. Through T ($T = 0, 1, 2, 3, \dots$) loops, we

employ an 1×1 convolution to fuse the output of each loop and obtain the final DIE feature vector. Subsequently, the DIE features are fed into the centering constrained cross-domain triplet loss (CCCDTL) function after L2-normalization. At the same time DIE features are forwarded to both identity classifier (IC) and domain classifier (DC) module.

The IC module consists of an FC layer and a Dropout layer [25]. This is a general multi-class classifier trained using standard cross-entropy loss function. This loss function is formulated as,

$$\begin{aligned} \mathcal{L}_{IC}(\mathbb{I}^s) = & -\frac{1}{|\mathbb{I}^s|} \sum_{I \in \mathbb{I}^s} (y_i \log \mathbb{P}(I) \\ & + (1 - y_i) \log(1 - \mathbb{P}(I))) \\ \text{with } & \mathbb{I}^s \cup \mathbb{I}^t = \mathbb{I} \end{aligned} \quad (1)$$

where \mathbb{I} represents images in a training mini-batch. \mathbb{I}^s denotes images from the source (labeled) domain and \mathbb{I}^t represents images from the target (unlabeled) domain. $\mathbb{P}(I)$ is the predicted probability of image I belonging to class y_i and $|\cdot|$ denotes the number of samples in set “.”.

B. Centering Constrained Cross-domain Triplet Loss

As triplet loss (TL) benefits from hard mining and metric learning, TL is a very common loss function in supervised Re-ID. In [8], Zhong et al. treat each two images from the target domain and the source domain as a negative pair, which enables TL to be used for cross-domain training. Following the assumption in [8], each image in target domain is assumed to have a different identities, since labels of the target domain are not available. The aforementioned strategy leads to a mistake when applied to cross-domain Re-ID. Even if two images from the target domain belong to the same person, they will be pushed away if TL is used in this way, which is demonstrated in Fig.2 a).

So, in order to alleviate this issue, we introduce a correction. More specifically, we mine hard positive pair only in the source domain and hard negative pair in both the source and the target domain. Meanwhile, we also introduce the Maximum Mean Discrepancy (MMD) distance to constrain target images which are pushed far away from the position where they should be. Finally, we propose a centering con-

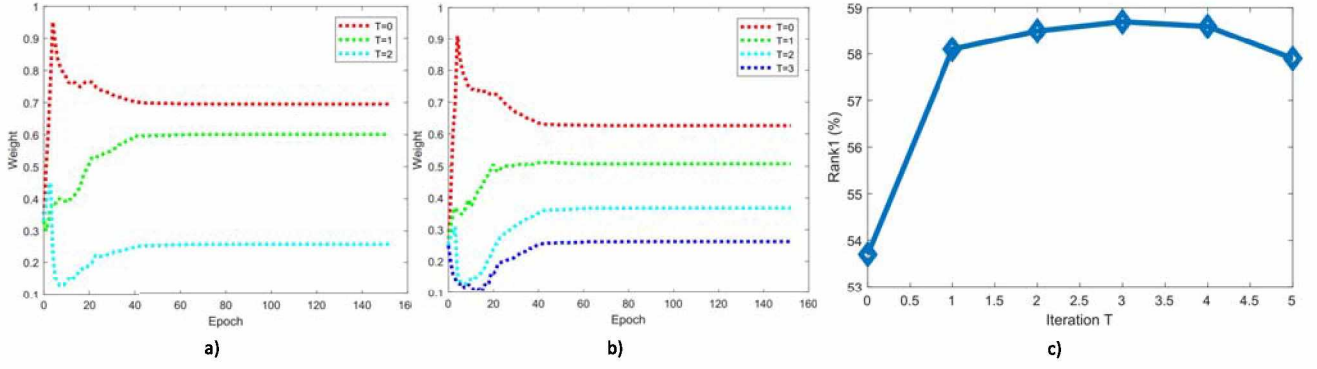


Fig. 3. a) The learned weight of fusing DIE feature ($T = 2$). b) The learned weight of fusing DIE feature ($T = 3$). c) Rank-1 accuracy of our model trained on the Duke dataset and tested on the Market dataset with different T s.

strained cross-domain triplet loss (CCCDTL) function to further improve discernment of embedding features by closing the farthest intra-class distance, pushing closest inter-class distance and minimizing the distance between the source and target distributions simultaneously, which is shown in Fig.2 b) and is formulated as,

$$\begin{aligned} \mathcal{L}_{CCCDTL}(\mathbb{I}) = & \sum_{I_a, I_p \in \mathbb{I}^s, I_n \in \mathbb{I}} \max\{D(\phi(I_a), \phi(I_p)) \\ & - D(\phi(I_a), \phi(I_n)) + m, 0\} \\ & + \lambda \times D\left(\frac{1}{|\mathbb{I}^s|} \sum_{I \in \mathbb{I}^s} \phi(I), \frac{1}{|\mathbb{I}^t|} \sum_{I \in \mathbb{I}^t} \phi(I)\right) \end{aligned} \quad (2)$$

with $\mathbb{I}^s \cup \mathbb{I}^t = \mathbb{I}$,

where ϕ is a complex non-linear map implemented by backbone network, which maps image sample to embedding. λ is a hyperparameter to balance the importance of two terms. I_a is an anchor point. I_p is the hardest (farthest) sample in the same class with I_a , and I_n is the hardest (closest) sample with a different class for I_a . m is a margin parameter and D is the Euclidean distance between two embedding feature vectors.

C. Domain-invariant Embedding by Gradient Reversal Layer

Inspired by unsupervised domain adaptive methods, we utilize the gradient reversal layer (GRL) in [17] to construct a domain classifier (DC) module to improve the domain-invariant capability of DIE features. Under the covariate shift assumption [26], we assume that there exist two distributions $S(I; L)$ and $T(I; L)$ (I and L denote images and labels respectively), which will be referred to as the source distribution and the target distribution. Both distributions are assumed to be very complex and unknown, and furthermore similar but different. In order to obtain a similar Re-ID accuracy on the target domain (unlabeled) as on the source domain (labeled), we should make the distributions S and T be similar. Unfortunately the distributions are unknown and can be very complex, which makes this problem difficult to solve. So, we reversely consider this problem that making two distributions as different as possible is equivalent to classifying them. With the help of gradient reversal layer

(GRL), we can transfer the classified supervised signal to an indiscriminate (domain-invariant) supervised signal, which is formulated as,

$$\begin{aligned} \mathcal{L}_{DC}(\mathbb{I}) = & -\frac{1}{|\mathbb{I}|} \sum_{I \in \mathbb{I}} (\Gamma_I \log \mathbb{P}(I) \\ & + (1 - \Gamma_I) \log(1 - \mathbb{P}(I))) \end{aligned} \quad (3)$$

$$\Gamma_I = \begin{cases} 1 & , I \in \mathbb{I}^t \\ 0 & , I \in \mathbb{I}^s \end{cases}$$

where Γ_I is an indicator function to index which domain image I belongs to. During the backpropagation processing, GRL makes the gradient negative and feeds it back to next layer, which is formulated as,

$$\theta \leftarrow \theta - \alpha \left(\frac{\partial \mathcal{L}_{IC}}{\partial \theta} + \frac{\partial \mathcal{L}_{CCCDTL}}{\partial \theta} - \frac{\partial \mathcal{L}_{DC}}{\partial \theta} \right), \quad (4)$$

where θ are all the parameters of the whole neural network and α is the step size of SGD.

D. Recurrent Top-down Attention

Most attention-based Re-ID models implement attention mechanism by utilizing extra neural network modules to predict where the model should focus on. These modules usually consist of Multi Layer Perceptrons (MLPs) and rely on the outputs of lower layers. The feature extraction and the attention prediction in the current layer work independently and at the same time. Intuitively, when an object catches our attention, our brain makes a decision to focus on the discriminative region, which should be a top-down (from brain to visual system) procedure. Similarly, our models also need to use the high-level feedback signal to guide feature extraction instead of using low-level features. Thus, we mimic the human visual attention process when addressing the Re-ID problem from in a complicated image, taking a first glimpse and then rethink several times, to optimize attention.

When a mini-batch of input images pass through all the layers, instead of immediately generating DIE feature vectors, a feedback module is deployed to propagate the supervised information to the bottom layers and update the network. On the one hand, intuitively, when two images with different

TABLE I
ABLATION STUDIES OF DOMAIN-INVARIANT EMBEDDING NETWORK UNDER DIFFERENT CONFIGURATIONS.

configurations	Duke==>Market				Market==>Duke			
	Rank-1	Rank-5	Rank-10	mAP	Rank-1	Rank-5	Rank-10	mAP
Base	45.1	62.7	70.1	20.8	32.9	49.7	54.9	17.0
Base+TL[8]	49.9	67.7	74.8	23.9	37.0	52.4	59.3	21.1
Base+CCCDTL	51.7	68.2	75.0	24.9	38.3	54.1	60.9	22.1
Base+DC	51.2	67.7	73.7	24.2	37.9	53.5	60.1	21.8
Base+CCCDTL+DC	53.6	69.9	76.1	26.3	39.8	55.3	62.5	22.3
Base+CCCDTL+DC+RTDA (T = 1)	58.1	74.1	80.9	24.8	44.4	60.1	66.1	24.2
Base+CCCDTL+DC+RTDA (T = 2)	58.5	74.9	81.4	26.9	45.8	61.6	67.8	26.1
Base+CCCDTL+DC+RTDA (T = 3)	58.7	75.4	81.9	27.1	46.7	62.5	68.3	26.4
Base+CCCDTL+DC+RTDA (T = 4)	58.6	75.2	81.5	27.0	46.4	62.3	68.0	26.2
Base+CCCDTL+DC+RTDA (T = 5)	57.9	73.9	80.4	26.3	44.1	61.9	67.8	24.0

identities have similar DIE features, they are not easy to be distinguished. Instead of outputting the feature vector directly, a better way is to recurrently guide the previous layers based on the primary DIE feature (when T equals to 0), such that the bottom layers can be strengthened or weakened to produce more discriminative features specifically for those identities that are difficult to distinguish. Furthermore, through aforementioned loss function, the DIE feature from the top layer will be more domain-invariant which is often come from high-level information. Thus we propose a new Recurrent Top-Down Attention (RTDA) module and allow DIE network to use the high-level feedback information for feature extraction.

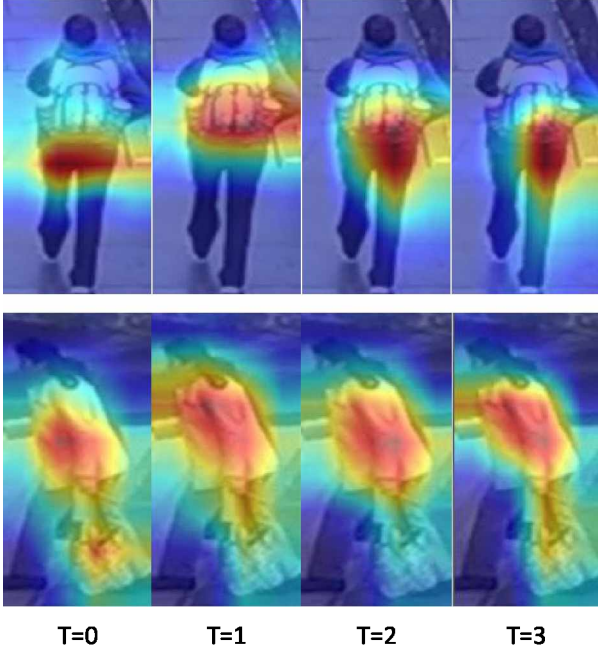


Fig. 4. Grad-CAM [27] visualization results for different Ts.

More specifically, we utilize the primary DIE feature during the first “glimpse” the image to predict the spatial positions of interest (spatial attention) and the weights of channels (channel attention) on the feature maps, and then make the network refocus on those regions and rethink em-

phasized or suppressed channels. In detail, spatial attention is implemented by three deconvolutional layers followed by a sigmoid layer, and channel attention consists of two convolutional layers and a sigmoid layer, as shown in Fig.1. We employ the Kronecker product to combine spatial- and channel-attention, which generates a spatial-channel mask with the same dimensions as the feature maps. After that these feature maps are updated by element-wise multiplication. After T times recurrent forward propagating, we fuse the DIE features from each loop by a weighted sum where the weights are learned by an 1×1 convolution and are initialized to $\frac{1}{T+1}$. Notably, our experiments on benchmark datasets clearly demonstrate the advantage of the RTDA algorithm in cross-domain Re-ID, which is reported in Section IV-B.

E. Multi-loss Joint Learning

In order to confirm that all modules work harmoniously and allow the proposed neural network to be trained in an end-to-end manner, we sum the three loss functions to form the final loss, which is written as follows:

$$\mathcal{L}_{final} = \mathcal{L}_{IC} + \beta_1 \mathcal{L}_{CCCDTL} + \beta_2 \mathcal{L}_{DC} \quad (5)$$

where β_1 and β_2 are hyper parameters to balance the importance of the three terms. Through cross validation, we set β_1 , β_2 and λ to 1, 0.1 and 1 respectively.

In addition, we adopt the stochastic gradient descent (SGD) method to update the parameters of the network while different learning rates are applied on different layers. More specifically, the weights of the pre-trained primary feature extractor should not be updated as fast as the other modules because we should keep the useful information acquired by training on ImageNet. Hence, we set the learning rate for the backbone network to a relatively small value, more specifically to 10^{-4} . For the other modules, IC, CCCDTL, DC and RTDA, the learning rate are 10^{-1} , 10^{-1} , 10^{-1} and 10^{-2} respectively.

IV. EXPERIMENTS AND ANALYSIS

A. Datasets

The performance of our proposed method on the task of Re-ID is evaluated on two popular benchmark datasets: Market-1501 [14] and DukeMTMC-reID [15].

TABLE II
COMPARISON WITH STATE-OF-THE-ART METHODS.

Methods	Duke==>Market				Market==>Duke			
	Rank-1	Rank-5	Rank-10	mAP	Rank-1	Rank-5	Rank-10	mAP
UMDL [28]	34.5	52.6	59.6	12.4	18.5	31.4	37.6	7.3
PUL [9]	45.5	60.7	66.7	20.5	30.0	43.4	48.5	16.4
SPGAN [10]	57.7	75.8	82.4	26.7	46.4	62.3	68.0	26.2
TJ-AIDL[18]	58.2	74.8	81.1	26.5	44.3	59.6	65.0	23.0
ours (T = 3)	58.7	75.4	81.9	27.1	46.7	62.5	68.3	26.4

Market-1501 consists of 32,668 labeled images of 1,501 identities collected from 6 camera views. All of the identities are divided into two parts: 12,936 images from 751 identities for training and 19,732 images from 750 identities for testing. During testing, 3368 query images from 750 identities are treated as probe for matching persons in the gallery.

DukeMTMC-reID is also a large-scale Re-ID dataset. It is collected from 8 cameras and contains 36,411 labeled images belonging to 1,404 identities. It consists of 16,522 training images from 702 identities, 2,228 query images from the other 702 identities, and 17,661 gallery images.

We use rank-1 accuracy and mean average precision (mAP) for our evaluation on both datasets. In the experiments, there are two source-target settings:

1. Target: Market-1501 / Source: DukeMTMC-reID.
2. Target: DukeMTMC-reID / Source: Market-1501.

B. Ablation Studie

In order to analyze the effectiveness of the proposed Domain-Invariant Embedding Network, we compare the baseline model with ten different configurations. For our model we use DIE features from different the Ts after L2-normalization as retrieved vectors in testing. The result of each experiment is reported in the each row of Table I.

Effectiveness of CCCDTL and DC module: As for the first configuration in Table I, our baseline model consists of the backbone network and the IC module, which is trained on the source datasets and directly evaluated on the test set of the target dataset. The second and third experiments aim at comparing the triplet loss function in [8] with our CCCDTL model. The results show that our method addresses the aforementioned mistake and shows a better performance. Furthermore, by adding an independent DC module into the base model, Rank-1 accuracy is increased by 6.1%. From Fig. 2 c) and d), we can see that the two distributions of the different domains blend into each other due to the effectivity of the DC module. Furthermore, our experiments show that all of the proposed modules are effectively increasing accuracy. Furthermore, combining with third, fourth and fifth row of Table I, the proposed modules do not conflict to each other, combining all loss functions achieves 53.6% at Rank-1 accuracy.

Effectiveness of RTDA module: Firstly, in order to investigate the influence of different $T = 1s$, we conducted several experiments using T ($T = 1, 2, 3, 4$ and 5). From Fig.3, it is obvious that increasing the number of recursive steps for information feedback allows the bottom layers to

receive richer top-down information. We observe from our experiments that after $T > 3$ the performance decreases since the model is overfitted. Empirically, we set $T = 3$ in the training phase of DIEN to compare with the state-of-the-art. Furthermore, combining Fig.3 a), b) and Fig.4 leads us to think that our model extracts coarse features which contain a large proportion of the information when $T = 0$. With the a increase in iterations, extracted features are more fine-grained with smaller proportions of information. Benefiting from aggregating multi-step features, the Re-ID performance of our model significantly increases again. Furthermore, from the Grad-CAM visualization results with different Ts in Fig.4, we can observe that domain-invariant features pay more attention on discriminative cues but not on the complete foreground. Finally, our proposed model achieves 58.7% at Rank-1 accuracy on the Market-1501 dataset. Hence, adding RTDA modules does help perform representation learning and it is cooperating with the DC and CCCDTL modules.

C. Comparison with State-of-the-art Methods

We compared our method with the state-of-the-art unsupervised learning methods. Table II presents the comparison when Market-1501 is the source set and Duke is the target set and viceversa. We compared with four unsupervised methods, including UMDL [28], PUL [9], SPGAN [10] and TJ-AIDL[18].

UMDL employed hand-crafted features and a multi-task dictionary learning method to learn cross-dataset feature, PUL is a typical post-processing method by reselecting training samples for fine-tuning a CNN model, SPGAN is a famous GAN-based baseline method for Person Re-ID, and TJ-AIDL is recently published and achieves the state-of-the-art result. Compared to the PUL method, our method achieves +13.2% higher rank-1 accuracy and a +6.6% improvement for mAP. As for the comparison with SPGAN, our method has +1.0% higher rank-1 accuracy and +0.2% higher mAP, while it is noted that GAN-based methods have significantly greater computational costs and memory consumption than our method and rely heavily on data augmentation. We also compare to TJ-AIDL and our results are slightly better than it, since TJ-AIDL given extra supervised information in the form of attributes of a person in the source dataset such as backpack or handbag et al. Our method without data augmentation has similar or better performance than all selected competitors.

V. CONCLUSION

In this paper, we proposed an end-to-end deep model, the Domain-invariant Embedding Network (DIEN), for solving cross-domain Re-ID tasks. Our DIEN utilizes both source-domain (labeled) datasets and target-domain (unlabeled) datasets as training data to explore the common cues of cross-domain Re-ID by jointly optimizing multiple loss functions. We also introduced a Recurrent Top-Down Attention module to refine the DIE features. Benefiting from the recurrent iteration, the model is able to extract more discriminative low-level features with the guidance from high-level information. With this proposed DIEN, we conducted experiments on internationally well known Market-1501 and DukeMTMC-reID datasets, and evaluated the effectiveness of our model in different configurations. Finally, compared to several recent unsupervised person Re-ID methods, the proposed DIEN achieved state-of-the-art performance and reduced the gap between supervised and unsupervised methods.

REFERENCES

- [1] Yumin Suh, Jingdong Wang, Siyu Tang, Tao Mei, and Kyoung Mu Lee. Part-aligned bilinear representations for person re-identification. In *European Conference on Computer Vision (ECCV)*, 2018.
- [2] Lingxiao He, Jian Liang, Haiqing Li, and Zhenan Sun. Deep spatial feature reconstruction for partial person re-identification: Alignment-free approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7073–7082, 2018.
- [3] Dapeng Chen, Dan Xu, Hongsheng Li, Nicu Sebe, and Xiaogang Wang. Group consistent similarity learning via deep crf for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8649–8658, 2018.
- [4] Jing Xu, Rui Zhao, Feng Zhu, Huaming Wang, and Wanli Ouyang. Attention-aware compositional network for person re-identification. 2018.
- [5] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016.
- [6] Weijian Deng, Liang Zheng, Guoliang Kang, Yi Yang, Qixiang Ye, and Jianbin Jiao. Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person reidentification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 6, 2018.
- [7] Zhun Zhong, Liang Zheng, Zhedong Zheng, Shaozi Li, and Yi Yang. Camera style adaptation for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5157–5166, 2018.
- [8] Zhun Zhong, Liang Zheng, Shaozi Li, and Yi Yang. Generalizing a person retrieval model hetero-and homogeneously. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 172–188, 2018.
- [9] Hehe Fan, Liang Zheng, Chenggang Yan, and Yi Yang. Unsupervised person re-identification: Clustering and fine-tuning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(4):83, 2018.
- [10] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. Domain separation networks. In *Advances in Neural Information Processing Systems*, pages 343–351, 2016.
- [11] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [12] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. *arXiv preprint arXiv:1711.08565*, 2017.
- [13] Vishal M Patel, Raghuraman Gopalan, Ruonan Li, and Rama Chellappa. Visual domain adaptation: A survey of recent advances. *IEEE signal processing magazine*, 32(3):53–69, 2015.
- [14] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1116–1124, 2015.
- [15] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. *arXiv preprint arXiv:1701.07717*, 3, 2017.
- [16] Minxian Li, Xiatian Zhu, and Shaogang Gong. Unsupervised person re-identification by deep learning tracklet association. *arXiv preprint arXiv:1809.02874*, 2018.
- [17] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. *arXiv preprint arXiv:1409.7495*, 2014.
- [18] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li. Transferable joint attribute-identity deep learning for unsupervised person re-identification. *arXiv preprint arXiv:1803.09786*, 2018.
- [19] Maurizio Corbetta and Gordon L Shulman. Control of goal-directed and stimulus-driven attention in the brain. *Nature reviews neuroscience*, 3(3):201, 2002.
- [20] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *CVPR*, volume 1, page 2, 2018.
- [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [22] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pages 1026–1034, 2015.
- [24] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention

- module. In *Proc. of European Conf. on Computer Vision (ECCV)*, 2018.
- [25] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from over-fitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [26] Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference*, 90(2):227–244, 2000.
- [27] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 618–626, 2017.
- [28] Peixi Peng, Tao Xiang, Yaowei Wang, Massimiliano Pontil, Shaogang Gong, Tiejun Huang, and Yonghong Tian. Unsupervised cross-dataset transfer learning for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1306–1315, 2016.