

Unsupervised Domain Adaptation in the Dissimilarity Space for Person Re-identification

Djebril Mekhazni, Amran Bhuiyan, George Ekladios, and Eric Granger

LIVIA, Dept. of Systems Engineering
École de technologie supérieure, Montreal, Canada
{djebril.mekhazni, amran.apece}@gmail.com, {george.ekladious,
eric.granger}@etsmtl.ca

Abstract. Person re-identification (ReID) remains a challenging task in many real-world video analytics and surveillance applications, even though state-of-the-art accuracy has improved considerably with the advent of deep learning (DL) models trained on large image datasets. Given the shift in distributions that typically occurs between video data captured from the source and target domains, and absence of labeled data from the target domain, it is difficult to adapt a DL model for accurate recognition of target data. DL models for unsupervised domain adaptation (UDA) are commonly designed in the feature representation space. We argue that for pair-wise matchers that rely on metric learning, e.g., Siamese networks for person ReID, the UDA objective should consist in aligning pair-wise dissimilarity between domains, rather than aligning feature representations. Moreover, dissimilarity representations are more suitable for designing open-set ReID systems, where identities differ in the source and target domains. In this paper, we propose a novel Dissimilarity-based Maximum Mean Discrepancy (D-MMD) loss for aligning pair-wise distances that can be optimized via gradient descent using relatively small batch sizes. From a person ReID perspective, the evaluation of D-MMD loss is straightforward since the tracklet information (provided by a person tracker) allows to label a distance vector as being either within-class (within-tracklet) or between-class (between-tracklet). This allows approximating the underlying distribution of target pair-wise distances for D-MMD loss optimization, and accordingly align source and target distance distributions. Empirical results with three challenging benchmark datasets show that the proposed D-MMD loss decreases as source and domain distributions become more similar. Extensive experimental evaluation also indicates that UDA methods that rely on the D-MMD loss can significantly outperform baseline and state-of-the-art UDA methods for person ReID. The dissimilarity space transformation allows to design reliable pair-wise matchers, without the common requirement for data augmentation and/or complex networks. Code is available on GitHub link: <https://github.com/djidje/D-MMD>

Keywords: Deep Learning, Domain Adaptation, Maximum Mean Discrepancy, Dissimilarity Space, Person Re-identification.

1 Introduction

Person re-identification (ReID) refers to the task of determining if a person of interest captured using a camera has the same identity as one of the candidates in the gallery, captured over different non-overlapping camera viewpoints. It is a key task in object recognition, drawing significant attention due to its wide range of applications, from video surveillance to sport analytics.

Despite the recent advances of ReID with DL models [10,17,20,24,26], and the availability of large amounts of labeled training data, person ReID still remains a challenging task due to the non-rigid structure of the human body, the different perspectives with which a pedestrian can be observed, the variability of capture conditions (e.g., illumination, blur), occlusions and background clutter. In practical video surveillance scenarios, the uncontrolled capture conditions and distributed camera viewpoints can lead to considerable intra-class variation, and to high inter-class similarity. The distribution of image data captured with different cameras and conditions may therefore differ considerably, a problem known in the literature as domain shift [18,28]. Given this domain shift, state-of-the-art DL models that undergo supervised training with a labeled image dataset (from the source domain) often generalize poorly for images captured in a target operational domain, leading to a decline in ReID accuracy.

Unsupervised domain adaptation (UDA) seeks resolve the domain shift problem by leveraging unlabeled data from the target domain (e.g., collected during a calibration process), in conjunction with labeled source domain data, to bridge the gap between the different domains. UDA techniques rely on different approaches, ranging from the optimization of a statistical criterion to the integration of an adversarial network, in order to learn robust domain-invariant representations from source and target domain data. Recently, several UDA methods have been proposed for pair-wise similarity matchers, as found in person ReID [5,14,27,29,30,33,34]. Common UDA approaches for metric learning employ (1) clustering algorithms for pseudo-labeling of the target data in the feature space, or (2) aligning feature representations of source and target data (either by minimizing some domain discrepancy or adversarial loss) [28]. These feature-based approaches are suitable for closed-set application scenarios, where the source and target domains share the same label space. However, this is not the case in open-set scenarios, where real-world person ReID systems are applied.

In this paper, we present a new concept for designing UDA methods that are suitable for pair-wise similarity matching in open-set person ReID scenarios. Instead of adapting the source model to unlabeled target samples in the feature representation space, UDA is performed in the dissimilarity representation space. As opposed to the common feature space, where a dimension represents a feature value extracted from one sample (i.e., a vector represents this sample measured over all features), the dissimilarity space consists of dissimilarity coordinates where each dimension represents the difference between two samples measured for a specific feature (i.e., a vector represents the Euclidean distance between two samples). Accordingly, the multiple clusters that represent different classes (i.e., ReID identities) in the feature representation space, are transformed to only two

clusters that represent the pair-wise within- and between-class distances. This transformation is more suitable for open-set ReID problems, when identities differ between the source and target domains, since the new label space has only two labels – pair-wise similar or dissimilar. Aligning the pair-wise distance distributions of the source and target domains in the dissimilarity space results in a domain-invariant pair-wise matcher.

The dissimilarity representation concept was recently introduced in [6], where a pseudo-labeling approach was proposed for UDA in still-to-video face recognition. This approach provided descent UDA results for problems with a limited domain shift. As a specific realization of the proposed concept, this paper focuses on a discrepancy-bases approach for dissimilarity-based UDA, that can provide a high level of accuracy for challenging problems with significant domain shift, as in ReID applications. To this end, we propose a variant of the common Maximum Mean Discrepancy (MMD) loss that is tailored for the dissimilarity representation space. The new Dissimilarity-based MMD (D-MMD) loss exploits the structure of intra- and inter-class distributions to align the source and target data in the dissimilarity space. It leverages tracklet¹ information to approximate the pair-wise distance distribution of the target domain, and thus estimate a reliable D-MMD loss for alignment of source and target distance distributions.

This paper contributes a novel D-MMD loss for UDA of DL models for person ReID. This loss allows to learn a domain-invariant pair-wise dissimilarity space representation, and thereby bridge the gap between image data from source and target domains (see Fig. 1). An extensive experimental analysis on three benchmark datasets indicates that minimizing the proposed D-MMD loss allows to align the source and target data distributions, which substantially enhances the recognition accuracy across domains. It also allows for designing reliable pair-wise matchers across domains, without the traditional requirement for data augmentation and/or complex networks.

2 Unsupervised Domain Adaptation for ReID

UDA focuses on adapting a model such that it can generalize well on an unlabeled target domain data while using a labeled source domain dataset. DL models for UDA seek to learn discriminant and domain-invariant representations from source and target data. They are generally based on either adversarial-, discrepancy-, or reconstruction-based approaches [28]. UDA methods have received limited attention in ReID because of their weak performance on benchmarks datasets compared to their supervised counterparts. Relying on a large-amount of annotated image data, and leveraging the recent success of deep convolutional networks, supervised ReID approaches [1,4,17,20,23] have shown a significant performance improvement, but UDA performance drops drastically when tested on different datasets and large domain shifts. To deal with this issue,

¹ A tracklet correspond to a sequence of bounding boxes that are captured over time for a same person in a camera viewpoint, and obtained using a person tracker.

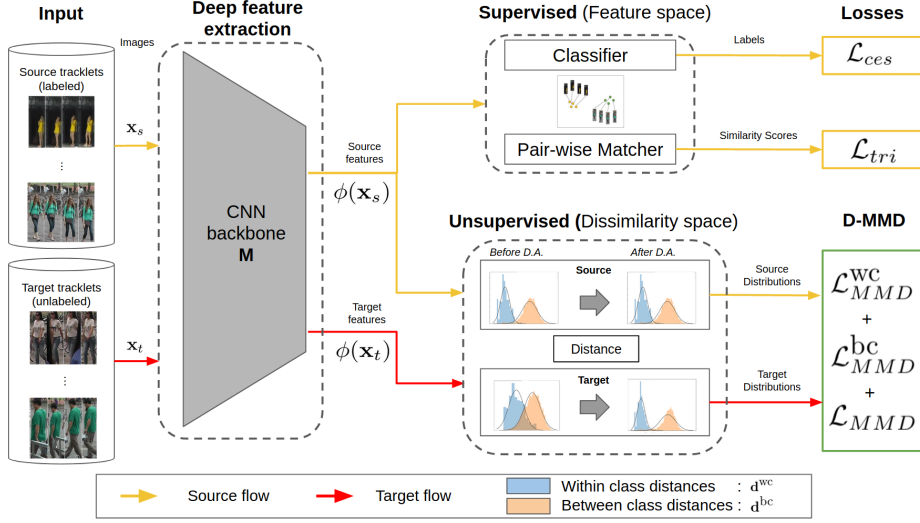


Fig. 1. Deep learning model for UDA using the proposed D-MMD loss. Labeled source images and unlabeled target images are input. First, the DL model for ReID undergoes supervised learning with source images. Upon reaching convergence, the backbone CNN can produce deep features from source and target images. Within-class (WC) and between-class (BC) dissimilarity distributions are produced for source and target domain data. Then, the D-MMD loss is applied between WC (resp. BC) source and WC (resp. BC) target. Supervised losses are also employed to ensure model stability.

representative methods use either clustering-based approach or domain-invariant feature learning based approach.

In clustering-based approaches [7,14], unlabeled target data are clustered to generate pseudo-labels, and then the network is optimized using the pseudo-labeled target data. Accordingly, performance of these approaches highly depend on the accuracy of clustering algorithms, and low accuracy can result in the propagation of noisy labels, and a corrupted model. In contrast, the domain-invariant feature learning based approaches [2,11,22,25,31] learn domain-invariant features. One approach is to define a discrepancy loss function that measures the domain shift in the feature space so that minimization of this loss decreases the domain shift, such in CORAL [22], MMD GAN [11], and WMMD [31]. Another approach for producing domain-invariant feature representations is through adversarial training, by penalizing a classifier’s ability to differentiate between source and target representations [2,25].

These approaches either employ pseudo-labeling using a specific set of labels (classes) that exist in the source domain, or represent samples of specific individuals similarly in both source and target domains. Therefore, these approaches are more suitable for closed-set application scenarios, where the source and target domains share the label space. Accordingly, these approaches can be

ineffective when applied to real-world person ReID applications that generally correspond to an open-set scenario. Indeed, individuals that appear in the target operational domain are typically different than those in design detests, or during the calibration phase.

To overcome these limitations of domain-invariant feature learning approaches, a different category of methods generate synthetic labeled data by transforming the source data to their style representative of target data [5,29,33,34]. However, performance of these approaches completely depends on the image generation quality. Other methods in the literature use labeled source data to train an initial deep ReID model, and then refine the trained model by clustering the target data [13,12,30]. These methods achieve a lower performance as they do not leverage the labeled source data to guide the adaptation procedure. Moreover, all aforementioned methods ignore the valuable knowledge that can be inferred from the underlying relations among target samples.

This paper addresses the limitations of the existing UDA methods for ReID through transferring the design space from the common feature representation space to the dissimilarity representation space, where open-set models can be easily adapted. This allows aligning the pair-wise distance distributions of the source and target domains. More specifically, this approach differs from the literature in two main aspects: (1) Unlike [11,22,31], we proposed to use D-MMD loss by exploiting the advantages of intra- and inter-class distributions along with global distributions. This allows dealing with the open-set application scenario exist in person ReID. (2) Our proposed approach does not rely on synthetic data augmentation as in [5,29,33,34], nor on the sensitivity of clustering algorithms as in [7,14].

3 Proposed Method

In this paper, a novel Dissimilarity-based Maximum Mean Discrepancy (D-MMD) loss is proposed for UDA of ReID systems. Rather than aligning source and target domains feature space, our D-MMD loss allows for the direct alignment of pair-wise distance distributions between domains. This involves jointly aligning the pair-wise distances from within-class distributions, as well as distances from between-class distributions. Both of these component contribute to accurate UDA for ReID systems based on a pair-wise similarity matcher, and have not been considered in other state-of-the-art methods. The proposed D-MMD loss allows to optimize pair-wise distances through gradient descent using relatively small batches.

Fig. 1 shows a DL model for UDA that relies on our D-MMD loss. For training, images $\mathbf{x}_s \in \mathbf{X}_s$ are sampled from the source domain \mathcal{D}_s , while images $\mathbf{x}_t \in \mathbf{X}_t$ are sampled from the target domain \mathcal{D}_t . During UDA, the CNN backbone model \mathcal{M} is adapted to produce a discriminant feature representation $\phi(\mathbf{x}_s)$ (resp. $\phi(\mathbf{x}_t)$) for input images, and the distances between input feature vectors allows estimating WC or BC distributions.

The underlying relations between source and target domain tracklets are employed to compute distributions of Euclidean distances based on samples of same identity (WC), \mathbf{d}^{wc} , and of different identities (BC), \mathbf{d}^{bc} . The D-MMD loss \mathcal{L}_{D-MMD} seeks to align the distance distributions of both domains through back propagation. The overall loss function \mathcal{L} for UDA is:

$$\mathcal{L} = \mathcal{L}_{\text{Supervised}} + \mathcal{L}_{D-MMD} \quad (1)$$

During inference, the ReID system performs pair-wise similarity matching. It is therefore relevant to optimize in the similarity space, and align target similarity distribution with well-separated intra/inter-class distribution from \mathcal{D}_s . The rest of this section provides additional details on the $\mathcal{L}_{\text{Supervised}}$ and \mathcal{L}_{D-MMD} loss functions.

3.1 Supervised Loss:

A model \mathcal{M} is trained through supervised learning on source data \mathbf{X}_s using a combination of a softmax cross-entropy loss with label smoothing regularizer (\mathcal{L}_{ces}) [24] and triplet loss (\mathcal{L}_{tri}) [10]. \mathcal{L}_{ces} is defined by Szegedy et al. [24] as:

$$\mathcal{L}_{\text{ces}} = (1 - \epsilon) \cdot \mathcal{L}_{\text{ce}} + \frac{\epsilon}{N}, \quad (2)$$

where N denotes total number of classes, and $\epsilon \in [0, 1]$ is a hyper-parameter that control the degree of label smoothing. \mathcal{L}_{ce} is defined as:

$$\mathcal{L}_{\text{ce}} = \frac{1}{K} \sum_i^K -\log \left(\frac{\exp(\mathbf{W}_{y_i}^T \mathbf{x}_i + b_{y_i})}{\sum_{j=1}^N \exp(\mathbf{W}_j^T \mathbf{x}_i + b_j)} \right) \quad (3)$$

where K is the batch size. Class label $y_i \in \{1, 2, \dots, N\}$ is associated with training image \mathbf{x}_i , the i^{th} training image. Weight vectors \mathbf{W}_{y_i} and bias b_{y_i} of last fully connected (FC) layer corresponds to class y of i^{th} image. \mathbf{W}_j and b_j are weights and bias of last FC corresponding of the j^{th} class ($j \in [1, N]$). \mathbf{W}_j and \mathbf{W}_{y_i} are respectively the j^{th} and y_i^{th} column of $\mathbf{W} = \{w_{ij} : i = 1, 2, \dots, F; j = 1, 2, \dots, N\}$, where F is the size of the last FC layer.

Triplet loss is also employed with hard positive/negative mining as proposed by Hermans et al. [10], where batches are formed by randomly selecting a person, and then sampling a number of images for each person. For each sample, the hardest positive and negative samples are used to compute the triplet loss:

$$\mathcal{L}_{\text{tri}} = \frac{1}{N_s} \sum_{\alpha=1}^{N_s} [m + \max(d(\phi(\mathbf{x}_\alpha^i), \phi(\mathbf{x}_p^i))) - \min_{i \neq j} (d(\phi(\mathbf{x}_\alpha^i), \phi(\mathbf{x}_n^j)))]_+ \quad (4)$$

where, $[\cdot]_+ = \max(\cdot, 0)$, m denotes a margin, N_s is the set of all hard triplets in the mini-batch, and d is the Euclidean distance. \mathbf{x}_j^i corresponds to the j^{th} image of the i^{th} person in a mini-batch. Subscript α indicates an anchor image, while p

and n indicate a positive and negative image with respect to that same specific anchor. $\phi(\mathbf{x})$ is the feature representation of an image \mathbf{x} .

For our supervised loss, we combine both the above losses.

$$\mathcal{L}_{\text{supervised}} = \mathcal{L}_{\text{ces}} + \lambda \cdot \mathcal{L}_{\text{tri}} \quad (5)$$

where λ is a hyper-parameter that weights the contribution of each loss term.

The softmax cross-entropy loss \mathcal{L}_{ces} defines the learning process as a classification task, where each input image is classified as one of the known identities in the training set. The triplet loss \mathcal{L}_{tri} allows to optimise an embedding where feature vectors are less similar different for inter-class images, and more similar intra-class images.

3.2 Dissimilarity-based Maximum Mean Discrepancy (D-MMD):

After training the model \mathcal{M} , we use it to extract feature representations from each source image $\mathbf{x}_s \in \mathbf{X}_s$, $\phi(\mathbf{x}_s)$, and target image $\mathbf{x}_t \in \mathbf{X}_t$, $\phi(\mathbf{x}_t)$. Then, the within-class distances, e.g., Euclidean or L_2 distances, between each different pair of images \mathbf{x}_i^u and \mathbf{x}_i^v of the same class i are computed:

$$d_i^{\text{wc}}(\mathbf{x}_i^u, \mathbf{x}_i^v) = \|\phi(\mathbf{x}_i^u) - \phi(\mathbf{x}_i^v)\|_2, \quad u \neq v \quad (6)$$

where $\phi(\cdot)$ is the backbone CNN feature extraction, and x_i^u is the image u of the class i . Similarly, the between-class distances are computed using each different pair of images \mathbf{x}_i^u and \mathbf{x}_j^z of the different class i and j :

$$d_{i,j}^{\text{bc}}(\mathbf{x}_i^u, \mathbf{x}_j^z) = \|\phi(\mathbf{x}_i^u) - \phi(\mathbf{x}_j^z)\|_2, \quad i \neq j \ \& \ u \neq z \quad (7)$$

Then, \mathbf{d}^{wc} and \mathbf{d}^{bc} are defined as the distributions of all distance values d_i^{wc} and $d_{i,j}^{\text{bc}}$, respectively, in the dissimilarity space.

The within-class (WC) and between-class (BC) distance samples of the source domain are computed directly using the source labels, so they capture the exact pair-wise distance distribution of the source domain. On the other hand, given the unlabeled target data, we leverage the tracklet information provided by a visual object tracker. We consider the frames within same tracklet as WC samples, and frames from different tracklets as BC samples. It is important to note that such tracklet information provide us with an approximation of the pair-wise distance distribution of the target domain since it lacks intra-class pairs from different tracklets or cameras.

Maximum Mean Discrepancy (MMD) [9] metric is used to compute the distance between two distribution:

$$\begin{aligned} \text{MMD}(P(A), Q(B)) &= \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n k(a_i, a_j) \\ &\quad + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m k(b_i, b_j) - \frac{2}{nm} \sum_{i=1}^n \sum_{j=1}^m k(a_i, b_j) \quad (8) \end{aligned}$$

where A (resp. B) is the source (resp. target) domain and $P(A)$ (resp. $Q(B)$) is the distribution of the source (resp. target) domain. $k(.,.)$ is a kernel (e.g. Gaussian) and a_i (resp. b_i) is sample i from A (resp. B). n and m are number of training examples from $P(A)$ and $Q(B)$, respectively.

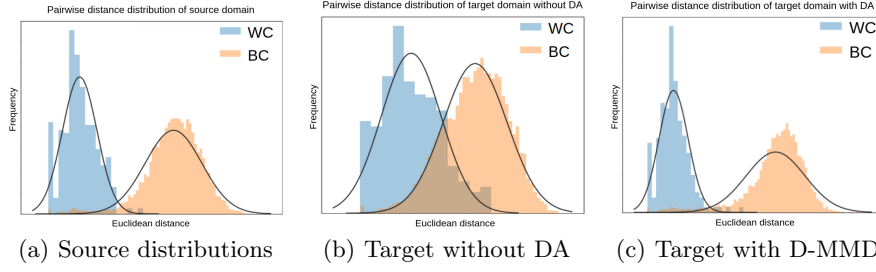


Fig. 2. Fig. 2(a) shows that the dissimilarity representation of the WC (blue) and BC (orange) distributions, where BC has larger Euclidean distances than WC because the model produces features closer for samples from same identities than that for samples from different identities. Fig. 2(b) shows a significant overlap when target data are represented using the initial source model, due to the intrinsic domain shift. Figure 2(c) shows that the target BC and WC distributions become aligned with the source distributions (Fig. 2(a)) after performing UDA.

To evaluate the divergence between two domains, MMD metric is applied to measure the difference from features produced by the source and target models are different using:

$$\mathcal{L}_{MMD} = MMD(\mathcal{S}, \mathcal{T}) \quad (9)$$

\mathcal{S} (\mathcal{T}) is defined as the distribution of the sources (target) images \mathbf{X}_s (\mathbf{X}_t) represented in the feature space.

Our method relies on the application of the MMD in the dissimilarity representation space instead of the common feature representation space. We define the \mathcal{L}_{MMD}^{wc} and \mathcal{L}_{MMD}^{bc} loss terms as follows:

$$\mathcal{L}_{MMD}^{wc} = MMD(\mathbf{d}_s^{wc}, \mathbf{d}_t^{wc}) \quad (10)$$

$$\mathcal{L}_{MMD}^{bc} = MMD(\mathbf{d}_s^{bc}, \mathbf{d}_t^{bc}) \quad (11)$$

Minimizing the above terms aligns the pair-wise distance distributions of the source and target domains, so that pair-wise distances from different domains are not deferential, and hence the source model works well in the target domain. Finally, our unsupervised loss function can be expressed as:

$$\mathcal{L}_{D-MMD} = \mathcal{L}_{MMD}^{wc} + \mathcal{L}_{MMD}^{bc} + \mathcal{L}_{MMD} \quad (12)$$

Algorithm 1 presents a UDA training strategy based on the D-MMD loss. Firstly, a supervised training phase runs for N^s epochs and produces a reference model \mathcal{M} using

Algorithm 1 UDA training strategy based on the D-MMD loss.

Require: labeled source data \mathbf{X}_s , and unlabeled target data \mathbf{X}_t
Load source data \mathbf{X}_s , and **Initialize** backbone model \mathcal{M}
for $l \in [1, N^s]$ epochs **do**
 for each mini-batch $B_s \subset \mathbf{X}_s$ **do**
 1) **Compute** \mathcal{L}_{ces} with Eq. 2, and \mathcal{L}_{tri} with Eq. 4
 2) **Optimize** \mathcal{M}
 end for
end for
Load target data \mathbf{X}_t , and **Load** backbone model \mathcal{M}
for $l \in [1, N^u]$ epochs **do**
 for each mini-batch $B_s \subset \mathbf{X}_s$ each mini-batch $B_t \subset \mathbf{X}_t$ **do**
 1) **Generate** \mathbf{d}_s^{wc} with Eq. 6, \mathbf{d}_s^{bc} with Eq. 7 and B_s
 2) **Generate** \mathbf{d}_t^{wc} with Eq. 6, \mathbf{d}_t^{bc} with Eq. 7 and B_t
 3) **Compute** $\mathcal{L}_{D-MMD}(B_s, B_t)$ with Eq. 12
 4) **Compute** $\mathcal{L}_{Supervised}$ with Eq. 5 using B_s
 5) **Optimize** \mathcal{M} based on overall \mathcal{L} (Eq. 1)
 end for
end for

the source domain data. Then, an unsupervised training phase runs for N^u epochs and aligns the target and source pair-wise distance distributions by minimizing the D-MMD loss terms defined by Eq. 9, Eq. 10, and Eq. 11. Note that the supervised loss $\mathcal{L}_{Supervised}$ is evaluated during domain adaptation to ensure the model \mathcal{M} remains aligned to a reliable source distribution \mathcal{S} over training iterations remaining.

Evaluating the D-MMD loss during the UDA training strategy involves computing the distances among each pairs of images. The computational complexity can be estimated as the number of within-class and between-class distance calculations. Assuming a common batch size of $|B|$ for training with source and target images and a number of occurrence of the same identity N_o , the total number of distance calculations is:

$$N_{\text{distances}} = N_{\text{distances}}^{wc} + N_{\text{distances}}^{bc} = (N_o - 1)! \frac{|B|}{N_o} + N_o \left(\frac{|B|}{N_o} - 1 \right)^2 \quad (13)$$

4 Results and Discussion

4.1 Experimental methodology:

For the experimental validation, we employ three challenging person ReID datasets, Market-1501 [32], DukeMTMC [21] and MSMT17 [29], and compare our proposed approach with state-of-the-art generative (GAN), tracklet-based, and domain adaptation methods for unsupervised person ReID.

Table. 1 describes three datasets for our experimental evaluation – Market1501, DukeMTMC and MSMT17. the **Market-1501** [32] dataset comprises labels generated using Discriminatively Trained Part-Based Models (DPM) [8]. It provides a realistic benchmark, using 6 different cameras, and around ten times more images than previously published datasets. The **DukeMTMC** [21] dataset is comprised of videos

Table 1. Properties of the three challenging datasets used in our experiments. They are listed according to their complexity (number of images, persons, cameras, and capture conditions, e.g., occlusions and illumination changes).

Datasets	# IDs	# cameras	# images	# train (IDs)	# gallery (IDs)	# query (IDs)	Annotation method	Crop size
Market-1501	1501	6	32217	12936 (751)	15913 (751)	3368 (751)	semi-automated (DPM)	128x64
Duke-MTMC	1812	8	36441	16522 (702)	17661 (1110)	2228 (702)	manual	variable
MSMT17	4101	15	126441	32621 (1041)	82161 (3060)	11659 (3060)	semi-automated (Faster R-CNN)	variable

captured outdoor at the Duke University campus from 8 cameras. **MSMT17** [29] is the largest and most challenging ReID dataset. It is comprised of indoor and outdoor scenarios, in the morning, noon and afternoon, and each video is captured over a long period of time.

For the supervised training, a Resnet50 architecture model is pretrained on ImageNet until convergence for both Hard-Batch Triplet and Softmax Cross-Entropy loss functions. Source domain videos are utilized for supervised training and evaluation. Then, the source and target training videos are used to perform UDA of the source model. Features are extracted from images of both domains using the Resnet50 CNN backbone (with a 2048 features vector size). To compute the BC and WC distributions, we randomly selected 4 occurrences of each class within batches of size 128. Given the nature of data, the tracklets are subject to greater diversity, with images from different viewpoints. The $D - MMD$ is then computed as described in Section 3, and backpropagation is performed using an Adam optimizer with a single step scheduler, decreasing the learning rate by 10 (initially 0.003) after every 20 epochs. In all steps, every image is resized to 256×128 before being processed.

Table 2 reports the upper bound accuracy for our datasets. To obtain this reference, we leveraged labeled source and target image data for supervised training. The ResNet50 model is initially trained using data from a first person ReID dataset (source domain), and then it is fine-tuned with training data from a second ReID datasets (target domain). Accuracy is computed with the target test sets of respective ReID datasets. We employed cross entropy loss with label smoothing regularizer 2 with $\epsilon = 0.1$, and triplet loss with a margin $m = 0.3$. To train DukeMTMC and Market1501, 30 epochs are required, while MSMT17 requires 59 epochs due to its larger-scale and complexity. Results in Table 2 confirms that MSMT17 is the most challenging dataset and shows lowest performance (63.2 % rank-1 accuracy).

Instead of optimizing the number of occurrences (frames) in a tracklet as a hyperparameter, we had to use a fixed number (4 occurrences) since the experimental datasets can sometime include only this number of frames per tracklet, and also for fair comparison with the SOA results. The metrics used for performance evaluation are the mean average precision (mAP), and rank-1, rank-5, rank-10 accuracy from the Cumulative Match Curve (CMC).

Table 2. Upper bound accuracy obtained after training on source data, and then fine-tuning on target data. Accuracy is measured with target training data.

Dataset source \rightarrow target	Accuracy			
	rank-1	rank-5	rank-10	mAP
DukeMTMC \rightarrow Market1501	89.5	95.6	97.1	75.1
Market1501 \rightarrow DukeMTMC	79.3	89.3	92.0	62.7
DukeMTMC \rightarrow MSMT17	63.2	77.5	82.0	33.9

4.2 Ablation study:

Table. 3 shows the impact on accuracy of the different loss terms. It is clear that \mathcal{L}_{MMD}^{bc} and \mathcal{L}_{MMD}^{wc} provide important information, with a slight improvement for the between-class (BC) component. Moreover, results show that a combination of both losses produces better results than when each term is employed separately. Moreover, while the classic feature-based \mathcal{L}_{MMD} had insignificant impact when employed separately (as observed in [14]), it helps when combined with the other terms. This can be explained by the fact that \mathcal{L}_{MMD} suffers from ambiguous association while dealing with domain shift that exists in open-set scenarios. Nevertheless, when the domain gap decreases to a reasonable limit (with the help of the proposed dissimilarity-based loss terms \mathcal{L}_{MMD}^{bc} and \mathcal{L}_{MMD}^{wc}), the feature-based loss starts to contributing ReID accuracy.

From source DukeMTMC to target Market1501, the margin of improvement while considering only the WC component over the baseline are 9.7% for Rank-1 accuracy and 8.2% for mAP and for only BC component 15.7 % for Rank-1 accuracy and 12.1 % for mAP. From source Market1501 to target DukeMTMC, we reach for the WC component 6.6 % for Rank-1 accuracy and 4.3% for mAP improvement compared to the baseline when for BC component we obtain 21.9% for Rank-1 accuracy and 16.9 % for mAP more than the baseline.

Table. 3 shows that a model adapted using only BC information is capable to produce better representation and leads to better results (51.8% Rank-1 accuracy) than when using only WC (45.8% Rank-1 accuracy) for the DukeMTMC to Market1501 transfer problem. In general, combining the different terms provides better results than when individual losses are employed.

Table 3. Ablation Study. Impact on accuracy of individual loss terms when transferring between the DukeMTMC and Market1501 domains. (The lower bound accuracy refers is obtained with the ResNet50 model trained on source data, and tested on target data, without domain adaptation.)

Setting	Loss Functions				Source: Duke Target: Market		Source: Market Target: Duke	
	\mathcal{L}_{sup}	\mathcal{L}_{MMD}^{wc}	\mathcal{L}_{MMD}^{bc}	\mathcal{L}_{MMD}	rank-1	mAP	rank-1	mAP
Lower Bound	✓	✗	✗	✗	36.1	16.1	23.7	12.3
A	✓	✓	✗	✗	45.8	24.3	30.3	16.6
B	✓	✗	✓	✗	51.8	28.4	45.6	29.2
C	✓	✓	✓	✗	66.6	45.4	60.5	42.9
D	✓	✓	✓	✓	70.6	48.8	63.5	46.0

In Section 3.2 (Fig.2), it is shown clearly that there is large overlap between the intra- and inter-class pair-wise distance distributions when using the initial source representation in the target domain. When applying the proposed method, the overlap significantly decreases and aligned with the source distributions. Fig 3 shows the reflection of such improvement of the pair-wise distance representation on the actual Re-ID problem. Before applying DA, there is much confusion between person representations which can be improved significantly with applying the proposed method.

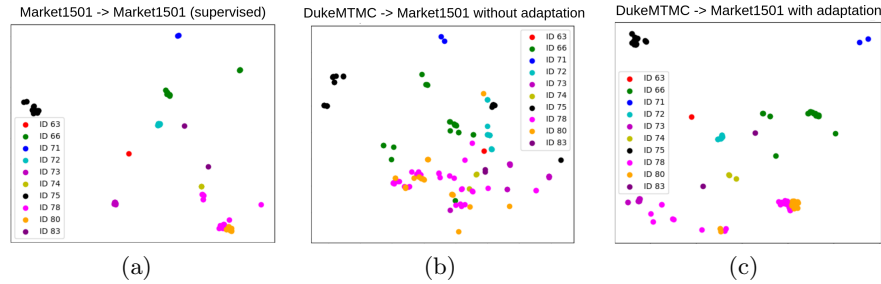


Fig. 3. T-SNE visualisations that show the impact of the original domain shift (3(a) versus 3(b)). Then 3(b) and 3(c) show the impact of employing our method to decrease domain shift, and accordingly improving the representation.

4.3 Comparison with state-of-art methods:

We compare our approach with state-of-the-art (SOTA) unsupervised methods on Market-1501, DukeMTMC-reID and MSMT17. Lower Bound refers to the domain shift without any adaptation. Table. 4 reports the comparison when tested on DukeMTMC and MSMT17 with Market1501 as the source, and Table. 4 reports results when DukeMTMC is the source.

PUL [7] and BUC [16] are clustering methods for pseudo-labeling of target data. Such approaches lead to poor performance. We outperform them by a large margin, 16.4% Rank-1 accuracy and 18.5% mAP more from Market1501 to DukeMTMC than BUC approach [16]. TAUDL [12] and UTAL [13] are two tracklet-based approaches for unsupervised person ReID. Due to their fully unsupervised behavior, they obtain worse results than our approach.

We also compare with other UDA approaches: PTGAN [29], SPGAN [5], ARN [15], TJ-AIDL [27] (attribute-based), [3] HHL [33], ECN [34], PDA-Net [14], Wu et al. [30], UDCA-CCE [19] (Camera-aware). Most of them are using data augmentation methods [34,33,29,5]. We are not using such techniques which are computationally expensive and require more memory. This also helps with problems that involve transferring from a small dataset to a larger and more complex dataset, which reflects a natural real-world application scenario. For the DukeMTMC to Market1501 transfer problem, we notice that DukeMTMC is a better initialization for simpler domains such Market1501, and it is easier to perform well in that sense (similar phenomenon for all other methods).

ECN and PDA-Net obtain better results on CMC metrics Rank-1 for this transfer problem DukeMTMC-Market1501 than ours (ECN [34] has 5.0% more and PDA-net

Table 4. ReID accuracy of the proposed and SOTA methods for UDA using Market1501 as source, and DukeMTMC and MSMT17 as targets. Accuracy is obtained on target datasets.

Methods	Source: Market1501								Conference or Journal
	DukeMTMC				MSMT17				
	r-1	r-5	r-10	mAP	r-1	r-5	r-10	mAP	
Lower Bound	23.7	38.8	44.7	12.3	6.1	12.0	15.6	2.0	–
PUL[7]	30.0	43.4	48.5	16.4	-	-	-	-	TOMM18
CFSM [3]	49.8	-	-	27.3	-	-	-	-	AAAI19
BUC [16]	47.4	62.6	68.4	27.5	-	-	-	-	AAAI19
ARN [15]	60.2	73.9	79.5	33.5	-	-	-	-	CVPR18-WS
UCDA-CCE [19]	47.7	-	-	31.0	-	-	-	-	ICCV19
PTGAN [29]	27.4	-	50.7	-	10.2	-	24.4	2.9	CVPR18
SPGAN+LMP[5]	46.4	62.3	68.0	26.2	-	-	-	-	CVPR18
HHL[33]	46.9	61.0	66.7	27.2	-	-	-	-	ECCV18
TAUDL[12]	61.7	-	-	43.5	28.4	-	-	12.5	ECCV18
UTAL[13]	62.3	-	-	44.6	31.4	-	-	13.1	TPAMI19
TJ-AIDL[27]	44.3	59.6	65.0	23.0	-	-	-	-	CVPR'18
Wu et al.[30]	51.5	66.7	71.7	30.5	-	-	-	-	ICCV19
ECN[34]	63.3	75.8	80.4	40.4	25.3	36.3	42.1	8.5	CVPR19
PDA-Net[14]	63.2	77.0	82.5	45.1	-	-	-	-	IEEE'19
D-MMD (Ours)	63.5	78.8	83.9	46.0	29.1	46.3	54.1	13.5	–

Table 5. ReID accuracy of the proposed and SOTA methods for UDA using DukeMTMC as source, and Market1501 and MSMT17 as targets. Accuracy is obtained on target datasets.

Methods	Source: DukeMTMC								Conference or Journal
	Market1501				MSMT17				
	r-1	r-5	r-10	mAP	r-1	r-5	r-10	mAP	
Lower Bound	36.6	54.5	62.9	16.1	11.3	20.6	25.7	3.7	–
PUL[7]	45.5	60.7	66.7	20.5	-	-	-	-	TOMM18
CFSM [3]	61.2	-	-	28.3	-	-	-	-	AAAI19
BUC [16]	66.2	79.6	84.5	38.3	-	-	-	-	AAAI19
ARN [15]	70.3	80.4	86.3	39.4	-	-	-	-	CVPR18-WS
UCDA-CCE [19]	60.4	-	-	30.9	-	-	-	-	ICCV19
PTGAN [29]	38.6	-	66.1	-	10.2	-	24.4	2.9	CVPR18
SPGAN+LMP [5]	57.7	75.8	82.4	26.7	-	-	-	-	CVPR18
HHL[33]	62.2	-	-	31.4	-	-	-	-	ECCV18
TAUDL[12]	63.7	-	-	41.2	28.4	-	-	12.5	ECCV18
UTAL[13]	69.2	-	-	46.2	31.4	-	-	13.1	TPAMI19
TJ-AIDL[27]	58.2	74.8	81.1	26.5	-	-	-	-	CVPR'18
Wu et al.[30]	64.7	80.2	85.6	35.6	-	-	-	-	ICCV19
ECN[34]	75.6	87.5	91.6	43.0	30.2	41.5	46.8	10.2	CVPR19
PDA-Net[14]	75.2	86.3	90.2	47.6	-	-	-	-	IEEE'19
D-MMD (Ours)	70.6	87.0	91.5	48.8	34.4	51.1	58.5	15.3	–

[14] has 4.6% higher Rank-1 accuracy for only 0.5% and 0.1% on Rank-5 and rank-10 accuracy regarding ECN). We outperform all other methods in mAP metrics (5.4% more than ECN). Since the D-MMD objective is to learn domain-invariant pair-wise dissimilarity representations, so success can be better measured using more global metrics, e.g., mAP, and this is validated by our results, where the proposed method produces best results using this metric. In contrast, CMC top-1 accuracy could be improved by training a pattern classifier (eg. MLP) to process the resulting distance vector. Nevertheless, when considering the opposite transfer problem, i.e., Market1501 to DukeMTMC, which is much more complex, our method provides best results for all metrics (0.3% on Rank-1, 3% on Rank-5, 3.5% on Rank-10 accuracy and 5.6% on mAP). We also outperform state-of-the-art methods on the most challenging dataset MSMT17 by 4.2% Rank-1, 9.6% Rank-5, 11.7% Rank-10 and 5.1% mAP when considering source DukeMTMC dataset. Similar results are observed using Market1501 as source.

Table 6. UDA accuracy of the proposed versus lower and upper bound approaches when transferring from MSMT17 (source) to Market1501 and DukeMTMC (targets).

Methods	Source: MSMT17							
	Market1501				DukeMTMC			
	rank-1	rank-5	rank-10	mAP	rank-1	rank-5	rank-10	mAP
Lower Bound	43.2	61.4	68.6	20.7	47.4	63.7	69.2	27.5
D-MMD (Ours)	72.8	88.1	92.3	50.8	68.8	82.6	87.1	51.6
Upper Bound	89.5	95.6	97.1	75.1	79.3	89.3	92.0	62.7

The proposed method can provide best performance for problems where the source domain consists in challenging data with high intra-class variability and high inter-class similarity (e.g., MSMT17) as compared to easier target domains (e.g. Market1501 and DukeMTMC). Such transfer problem is less explored in the literature, so in Table. 6 we compare our results with only the lower and upper bounds. With this setup (i.e. source is the most challenging dataset MSMT17) we obtained best results (better than these reported on Tables 4 and 5): 77.8% Rank-1 accuracy and 50.8% mAP for Market1501 and 68.8% Rank-1 accuracy and 51.6% mAP for DukeMTMC.

5 Conclusion

In this paper, we proposed a novel dissimilarity-based UDA approach for person ReID using MMD loss to reduce the gap between domains in the dissimilarity space. The core idea is to exploit the advantages of using within and between-class distances that effectively capture the underlying relations between domains which has never been explored in the state-of-the-art. To that end, we align the within- and between-class distance distributions for the source and target domains to produce effective Re-ID models for the target domain. Experiments on three challenging ReID datasets prove the effectiveness of this new approach as it outperforms state-of-the-art methods. Moreover, our proposed loss is general and can be applied to different feature extractors and applications.

References

1. Bhuiyan, A., Liu, Y., Siva, P., Javan, M., Ayed, I.B., Granger, E.: Pose guided fusion for person re-identification. In: WACV (2020)
2. Bousmalis, K., Trigeorgis, G., Silberman, N., Krishnan, D., Erhan, D.: Domain separation networks. In: NIPS (2016)
3. Chang, X., Yang, Y., Xiang, T., Hospedales, T.M.: Disjoint label space transfer learning with common factorised space. In: AAAI (2019)
4. Chen, B., Deng, W., Hu, J.: Mixed high-order attention network for person re-identification. In: ICCV (2019)
5. Deng, W., Zheng, L., Ye, Q., Kang, G., Yang, Y., Jiao, J.: Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification. In: CVPR (2018)
6. Ekladios, G., Lemoine, H., Granger, E., Kamali, K., Moudache, S.: Dual-triplet metric learning for unsupervised domain adaptation in video-based face recognition. In: IJCNN (2020)
7. Fan, H., Zheng, L., Yan, C., Yang, Y.: Unsupervised person re-identification: Clustering and fine-tuning. *ACM Trans. on Multimedia Computing, Communications, and Applications* **14**(4), 1–18 (2018)
8. Felzenszwalb, P.F., Girshick, R.B., McAllester, D., Ramanan, D.: Object detection with discriminatively trained part-based models. *IEEE Trans. on Pattern Analysis and Machine Intelligence* **32**(9), 1627–1645 (2009)
9. Gretton, A., Borgwardt, K.M., Rasch, M.J., Schölkopf, B., Smola, A.: A kernel two-sample test. *Journal of Machine Learning Research* **13**(1), 723–773 (2012)
10. Hermans, A., Beyer, L., Leibe, B.: In defense of the triplet loss for person re-identification (2017)
11. Li, C.L., Chang, W.C., Cheng, Y., Yang, Y., Póczos, B.: Mmd gan: Towards deeper understanding of moment matching network. In: NIPS (2017)
12. Li, M., Zhu, X., Gong, S.: Unsupervised person re-identification by deep learning tracklet association. In: ECCV (2018)
13. Li, M., Zhu, X., Gong, S.: Unsupervised tracklet person re-identification. *IEEE TPAMI* (2019)
14. Li, Y.J., Lin, C.S., Lin, Y.B., Wang, Y.C.F.: Cross-dataset person re-identification via unsupervised pose disentanglement and adaptation. In: ICCV (2019)
15. Li, Y.J., Yang, F.E., Liu, Y.C., Yeh, Y.Y., Du, X., Frank Wang, Y.C.: Adaptation and re-identification network: An unsupervised deep transfer learning approach to person re-identification. In: CVPR Workshops (2018)
16. Lin, Y., Dong, X., Zheng, L., Yan, Y., Yang, Y.: A bottom-up clustering approach to unsupervised person re-identification. In: AAAI (2019)
17. Luo, H., Gu, Y., Liao, X., Lai, S., Jiang, W.: Bag of tricks and a strong baseline for deep person re-identification. In: CVPR Workshops (2019)
18. Nguyen-Meidine, L.T., Granger, E., Kiran, M., Dolz, J., Blais-Morin, L.A.: Joint progressive knowledge distillation and unsupervised domain adaptation. In: IJCNN (2020)
19. Qi, L., Wang, L., Huo, J., Zhou, L., Shi, Y., Gao, Y.: A novel unsupervised camera-aware domain adaptation framework for person re-identification. In: ICCV (2019)
20. Quan, R., Dong, X., Wu, Y., Zhu, L., Yang, Y.: Auto-reid: Searching for a part-aware convnet for person re-identification. In: ICCV (2019)
21. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: ECCV (2016)

22. Sun, B., Saenko, K.: Deep coral: Correlation alignment for deep domain adaptation. In: ECCV (2016)
23. Sun, Y., Zheng, L., Yang, Y., Tian, Q., Wang, S.: Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In: ECCV (2018)
24. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. In: CVPR (2016)
25. Tzeng, E., Hoffman, J., Saenko, K., Darrell, T.: Adversarial discriminative domain adaptation. In: CVPR (2017)
26. Wang, G., Lai, J., Huang, P., Xie, X.: Spatial-temporal person re-identification. In: AAAI (2019)
27. Wang, J., Zhu, X., Gong, S., Li, W.: Transferable joint attribute-identity deep learning for unsupervised person re-identification. In: CVPR (2018)
28. Wang, M., Deng, W.: Deep visual domain adaptation: A survey. *Neurocomputing* **312**, 135–153 (2018)
29. Wei, L., Zhang, S., Gao, W., Tian, Q.: Person transfer gan to bridge domain gap for person re-identification. In: CVPR (2018)
30. Wu, A., Zheng, W.S., Lai, J.H.: Unsupervised person re-identification by camera-aware similarity consistency learning. In: ICCV (2019)
31. Yan, H., Ding, Y., Li, P., Wang, Q., Xu, Y., Zuo, W.: Mind the class weight bias: Weighted maximum mean discrepancy for unsupervised domain adaptation. In: CVPR (2017)
32. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: ICCV (2015)
33. Zhong, Z., Zheng, L., Li, S., Yang, Y.: Generalizing a person retrieval model hetero- and homogeneously. In: ECCV (2018)
34. Zhong, Z., Zheng, L., Luo, Z., Li, S., Yang, Y.: Invariance matters: Exemplar memory for domain adaptive person re-identification. In: CVPR (2019)