# Fully Unsupervised Person Re-identification via Selective Contrastive Learning

Bo Pang, Deming Zhai, *Member, IEEE,* Junjun Jiang, *Member, IEEE,* and Xianming Liu, *Member, IEEE*

arXiv:2010.07608v2 [cs.CV] 4 Mar 2021

*Abstract*—Person re-identification (ReID) aims at searching the same identity person among images captured by various cameras. Existing fully-supervised person ReID methods usually suffer from poor generalization capability caused by domain gaps. Unsupervised person ReID attracts a lot of attention recently, due to it works without intensive manual annotation and thus shows great potential of adapting to new conditions. Representation learning plays a critical role in unsupervised person ReID. In this work, we propose a novel selective contrastive learning framework for fully unsupervised feature learning. Specifically, different from traditional contrastive learning strategies, we propose to use multiple positives and adaptively selected negatives for defining the contrastive loss, enabling to learn a feature embedding model with stronger identity discriminative representation. Moreover, we propose to jointly leverage global and local features to construct three dynamic memory banks, among which the global and local ones are used for pairwise similarity computation and the mixture memory bank are used for contrastive loss definition. Experimental results demonstrate the superiority of our method in unsupervised person ReID compared with the state-of-the-arts.

*Index Terms*—Person Re-identification Unsupervised learning Contrastive learning

## I. INTRODUCTION

**P**ERSON re-identification (ReID), also referred to as person retrieval, aims at searching the same identity person among images captured by various cameras at different time and locations. Thanks to the rapid development of convolutional neural networks (CNN), in recent years, the performance of person ReID is improved remarkably by using discriminative features from labeled person images. However, the success of such systems relies on a large amount of labeled data, which is often prohibitively expensive to acquire. As a result, a large research effort is currently focused on unsupervised systems without leveraging intensive manual supervision, which attracts a lot of attention due to the great potential of adapting to new conditions.

This effort includes recent advances on transfer learning and unsupervised learning. Among them, cross-domain transfer learning, also called domain adaptation, offers an effective manner to reduce the labelling cost. The basic idea is to learn an attribute-semantic and identity discriminative representation from a labeled source dataset, which is then transferred to the target domain for person ReID [1, 2, 3, 4]. However, the performance of this approach largely relies on the assumption

B. Pang, D. Zhai, J. Jiang and X. Liu are with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin 150001, China, and also with the Peng Cheng Laboratory, Shenzhen 518052, China (cspb@hit.edu.cn; zhaideming@hit.edu.cn; jiangjunjun@hit.edu.cn; csxm@hit.edu.cn).

that there is sufficient knowledge overlap between the source and target domains, which is not always valid however. If the source domain exhibits significantly different characteristics with the target domain, the ReID performance would degrade heavily. Moreover, this kind of methods require a large amount of annotated source data, which are thus not purely unsupervised.

The fully unsupervised learning based approach receives more and more attention recently, since it works without leveraging any labeled data and thus shows better potential to deploy person ReID for real-world applications. The basic idea of this line is to alternate between predicting pseudo labels by clustering or classification and training the network with generated pseudo classes. For instance, Lin *et al.* [5] proposed a bottom-up clustering framework that iteratively trains a network based on the pseudo labels generated by unsupervised clustering. Wang *et al.* [6] formulated unsupervised person ReID as a multi-label classification task to progressively seek true labels. Lin *et al.* [7] proposed a classification network with softened labels to eliminate errors incurred from hard quantization in clustering. However, the performance of these methods relies on the accuracy of label prediction, which is a non-trivial task under unsupervised setting.

Instead of explicit label prediction, in this paper, we concentrate on contrastive self-supervised visual representation learning, taking advantage of the principle that a good feature representation model should map images of the same person closer to each other, while push images of different identities apart away. Specifically, we propose a novel selective contrastive learning framework with dynamic memory banks, which is specially tailored for the task of unsupervised person ReID. Although there is also some work attempts to tackle person ReID based on contrast learning such as [4], which is upon the domain adaption paradigm as thus requires a labeled source domain. In contrast, the proposed method is purely unsupervised.

The technical contributions of this work are three-fold:

- Considering that a person may be captured by various cameras, *i.e.*, each person may have multiple images in the training set, one key technical novelty is that we choose multiple positives for each anchor, as opposed to SimCLR [8], MoCo [9] and [4] that use only a single positive to define the contrastive loss. Moreover, different from the conventional contrastive learning strategies that take all samples except the positive as the negatives [4, 8, 9], we propose to select samples that are plausibly similar to the anchor as the negatives, so as to improve the discrimination ability of representation learning. More

specifically, using the defined distance metric, we rank the similarity order between the anchor and all training samples, according to which we divide the training set into three subsets: *similar set*, *borderline set*, and *dissimilar set*. By taking samples in *similar set* as the positives and samples in *borderline set* as the negatives, we define the contrastive loss to encourage the feature embedding function to produce closely aligned representations to all images of the same identity.

- The other contribution of this work is that we propose three dynamic dictionaries which jointly leverage global and local discriminative information for unsupervised representation learning. The global feature is widely used in existing unsupervised ReID, such as [5, 6]. However, it suffers from discriminative information loss in some cases, leading to images of different persons may have similar feature representations. On the other hand, the part-level features offer fine-grained discriminative information for pedestrian image description [10]. However, compared with the global feature, the local features bring much more search freedom, making the optimization of feature representation learning hard to converge. Moreover, learning discriminative local features requires that parts should be accurately located. It can be done either by external assistance from human pose estimation [11] or well-designed partition strategies [10, 12], which are expensive and hinder the generalization in practical applications. Considering their respective limitations, we propose to jointly leverage the global feature and the local features for defining distance metric and contrastive loss.

- By combining the above contributions, we propose an effective unsupervised person ReID algorithm. Our scheme achieves encouraging performance with respect to rank-1 and mAP so far on public Market-1501, DukeMTMC-reID, DukeMTMC-VideoReID and MARS. For instance, we achieve rank-1 accuracy of 82.2%, significantly outperforming the latest unsupervised person ReID methods SNR [13], SSLR [7], MMCL [6] and TSSL [14] by 15.5%, 10.5%, 15.6% and 11%, respectively.

The paper is organized as follows: Section 2 overviews some related works. Section 3 introduces the proposed scheme, including the representation learning framework, positives and negatives sampling strategy, the defined contrastive loss and the optimization strategy, and the memory bank update strategy. Section 4 provides the experimental results and ablation study. We conclude this paper in Section 5.

## II. RELATED WORK

### A. Unsupervised Domain Adaptation Person ReID

Transfer learning is a common strategy for addressing unsupervised person ReID. These domain adaption methods [1, 2, 3, 13, 15] attempt to tackle the unsupervised person ReID problem on the target unlabeled dataset by leveraging other dataset's labeled information. TJ-AIDL model [1] aims at learning an attribute-semantic and identity-discriminative feature representation space which is transferrable to the any unlabelled target dataset. HHL [15] aims to improve the

generalization ability of re-ID models on the target testing set with enforcing two properties, camera invariance and domain connectednes, simultaneously. Thanks to the development of the Generative Adversarial Network (GAN), this type of style transfer network is used for cross-domain transfer learning for unsupervised person ReID. SPGAN [2] generates transferred images from the source labeled dataset and then do the supervised learning with two constraints which are self-similarity of an image before and after translation and domain-dissimilarity of a translated source image and a target image. ATNet [3] decomposes the complicated cross-domain transfer into a set of factor-wise sub-transfers, each of which concentrates on style transfer with respect to a certain imaging factor, e.g., illumination, resolution and camera view etc. Considering poor generalization capability caused by domain gaps with existing methods, Baseline-SNR [13] filter out identity-irrelevant interference and learn domain-invariant person representations. In [4], a hybrid memory is proposed to encode all available information from both source and target domains for feature learning, which achieves the best unsupervised person ReID performance so far. Although achieving promising performance, these methods require an annotated source dataset. In contrast, our work focuses on purely unsupervised person ReID, which only relies on the unlabeled target dataset.

### B. One Shot Person ReID

Methods based on one-shot learning for person ReID attempt to solve the problem with condition where each identity has only one labeled example and many unlabeled examples. EUG [16] and ProLearn [17] propose to gradually and steadily improve the discriminative capability of the CNN via stepwise learning. Especially, for video-based person re-identification, RACE [18] firstly adopt anchor sequences to formulate an anchor graph. And then for accurately estimate labels from unlabeled sequences with noisy frames, robust anchor embedding is introduced based on the regularized affine hull. These methods solved the cost of annotation in a degree compared with the supervised methods and the domain adaptation based methods. But still, the labeled information is needed.

### C. Fully Unsupervised Person ReID

For fully unsupervised person ReID, traditional methods [19, 20] utilize hand-craft features, which are hardly designed to be discriminative by hand. Recently, the cluster based methods and the mutli-label based methods estimate pseudo labels to train the neural network. BUC [5] jointly optimize a convolutional neural network (CNN) and the relationship among the individual samples with bottom-up clustering procedure. TSSL [14] thought these bottom-up clustering methods merely utilise suboptimal global clustering. They design a comprehensive unsupervised learning objective that accounts for tracklet frame coherence, tracklet neighbourhood compactness, and tracklet cluster structure in a unified formulation, which is capable of capitalising directly from abundant unlabelled tracklet data. Wang *et al.* formulate the problem as a multi-label classification task to progressively seek true labels and adopt the memory-based multi-label classification

loss (MMCL) to boost the ReID model training efficiency. SSLR [7] proposed a similarity learning framework with softened labels to relief the hard quantization loss in clustering. However, these methods rely on the accuracy of the pseudo labels. Our proposed framework will adopt contrastive self-supervised visual representation learning, mapping the person images with same identity close to each other while push the person with different identity apart away.

### D. Unsupervised Representation Learning

Unsupervised representation learning utilize unlabeled data to learn an effective embedding space for downstream tasks like image classification and etc. Gidaris *et al.* [21] propose to learn image features by training ConvNets to recognize the 2d rotation that is applied to the image that it gets as input. Recently, contrastive learning has attracted a lot of attention. Such methods aim at mapping the representation close to the positives and away from the negatives. MoCo [9] build a large and consistent dictionary on-the-fly that facilitates contrastive unsupervised learning. SimCLR [8] propose a simple framework for contrastive learning of visual representations. Both the Moco and the SimCLR use only one positive to conduct contrastive loss. In our paper, we propose a new contrastive learning framework with selective positives and negatives and three dynamic dictionaries will be conducted to help the training process.

### III. METHODOLOGY

Given a training set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N\}$, the goal of unsupervised person ReID is to learn a model $\mathcal{F}(\boldsymbol{\theta}; \mathbf{x}_i)$ for visual feature representation without using any manual annotation, where parameters related to $\mathcal{F}$ are denoted as $\boldsymbol{\theta}$. The learned representation model is then applied on the query image $\mathbf{x}^q$ and the gallery set $\hat{\mathcal{X}} = \{\hat{\mathbf{x}}_1, \hat{\mathbf{x}}_2, ..., \hat{\mathbf{x}}_M\}$, so as to derive the query result by ranking distance between features of the query and all gallery images. It is clear to see that, unsupervised representation learning plays a central role in unsupervised person ReID.

The feature representation learning is challenging in unsupervised setting, wherein the critical problem is how to learn to discriminate between individual images, without any notion of semantic categories. A good feature representation model should map images of the same person closer to each other, while push images of different identities apart away. In this work, inspired by recent progress of self-supervised learning, we address the unsupervised feature learning problem from the perspective of contrastive learning.

### A. Representation Learning Framework

As illustrated in Fig. 1, our representation learning framework consists of the following components:

- *A neural network base encoder* $\mathbf{E}(\cdot)$, which is leveraged to extract feature vector for a given person image $\mathbf{x}_i$. It allows various choices of the network architecture without any constraints. Here we adopt ResNet-50 without the last classification layer as the backbone, whose

parameters are pre-trained on ImageNet, to obtain feature map $\mathbf{f}_i = \mathbf{E}(\mathbf{x}_i) \in \mathbb{R}^{2048}$.

- *Pooling operator* $\mathbb{AP}(\cdot)$, for which we consider both global and local average pooling to obtain global and local features. The global feature $\mathbf{f}_i^g$ is obtained by average pooling (AP) of feature map $\mathbf{f}_i$ of the whole image:

$$\mathbf{f}_i^g = \mathbb{AP}(\mathbf{f}_i) \tag{1}$$

However, the global feature suffers from discriminative information loss in some cases, leading to images of different identities may have similar feature representation. We further consider the part-level features, which offer fine-grained discriminative information for pedestrian image description [10]. We obtain part-level feature maps by equally partitioning the global feature map $\mathbf{f}_i$ into $N_l = 8$ horizontal stripes. The local features are then obtained by applying average pooling on each part:

$$\mathbf{f}_{i,j}^l = \mathbb{AP}(\mathbf{f}_{i,j}), j = 1, \cdots, N_l \tag{2}$$

- *A small projection neural network* $\mathbf{P}(\cdot)$, which is a learnable nonlinear operation that transforms image features to the latent space where the contrastive learning is conducted. $\mathbf{P}(\cdot)$ is defined with fully-connected (FC) layer, batch normalization (BN) layer and L2-normalization (L2-Norm) layer. The usage of $\mathbf{P}(\cdot)$ is demonstrated to be beneficial to define the contrastive loss [8].

- *Global and local memory bank* $\mathcal{M}^g$ and $\mathcal{M}^l$, which are leveraged to store global and local features for pairwise similarity computation. The keys in global memory bank $\mathcal{M}^g$ are defined as:

$$\mathbf{v}_i^g = \mathbf{P}\left(\mathbb{AP}(\mathbf{f}_i)\right), i = 1, \cdots, N \tag{3}$$

And the keys in local memory bank $\mathcal{M}^l$ are defined as:

$$\mathbf{v}_{i,j}^l = \mathbf{P}\left(\mathbb{AP}(\mathbf{f}_{i,j})\right), i = 1, \cdots, N; j = 1, \cdots, N_l \tag{4}$$

The global and local memory banks are then used to define global and local distance metrics respectively, which are further coupled with cross-camera encouragement term [7] to define the total distance metric for pairwise similarity computation, as shown in Eq. (11). According to the ranked similarity order, we identify the positives $\mathcal{K}_+$ and the negatives $\mathcal{K}_-$ that are used to define the contrastive loss. The details will be elaborated in the following subsection.

- *Mixture memory bank* $\mathcal{M}^t$, which is required to store features of training images to define the contrastive loss, so as to maximize the similarity between representations of the anchor and the positives, while minimizing that of the anchor and the negatives. To improve the discrimination ability, we construct a mixture memory bank $\mathcal{M}^t$ which includes the fusion of the global and local features as keys. It is worth noting that, for the local feature in $\mathcal{M}^t$, instead of using the one defined in Eq. (4), we turn to concatenate all $N_l$ ones to form a single local feature:

$$\mathbf{v}_i^l = \mathbb{CONCAT}\left(\{\mathbb{AP}(\mathbf{f}_{i,j})\}_{j=1}^{N_l}\right) \tag{5}$$
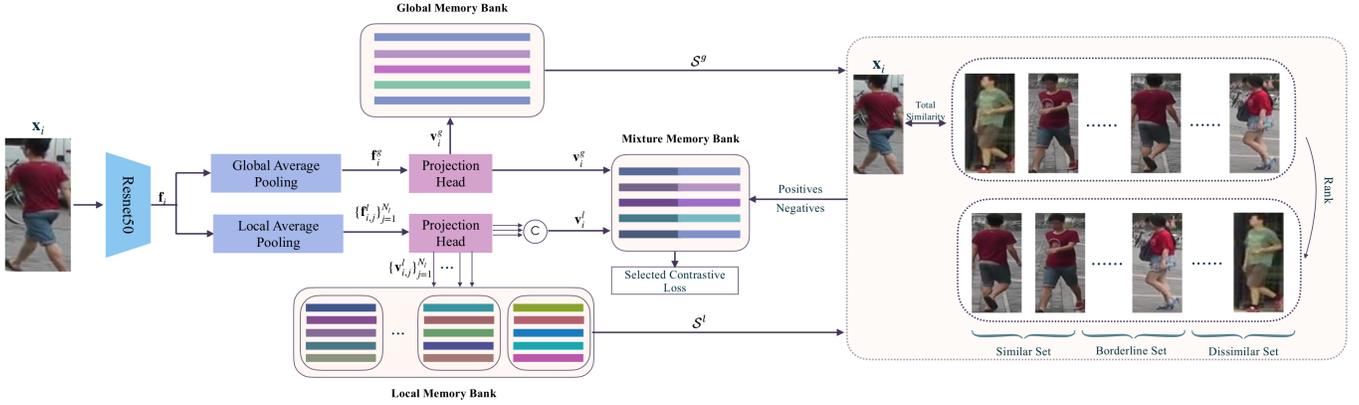
Fig. 1: The overall framework of our proposed unsupervised person ReID.

The mixture memory bank $\mathcal{M}^t$ is initialized with all zeros. We update its keys corresponding to all positives by fusing with the global and local features of the anchor $\mathbf{x}_i$ progressively. Specifically, the positive keys are firstly updated with global feature $\mathbf{v}_i^g$:

$$\mathcal{M}^t[k] = \|\frac{\mathcal{M}^t[k] + \mathbf{v}_i^g}{2}\|_2, k \in \mathcal{K}_+ \tag{6}$$

which are further updated with local feature $\mathbf{v}_i^l$:

$$\mathcal{M}^t[k] = \|\frac{\mathcal{M}^t[k] + \mathbf{v}_i^l}{2}\|_2, k \in \mathcal{K}_+ \tag{7}$$

where $\| \cdot \|_2$ represents L2-normalization. In this way, the mixture memory bank jointly employs the global and local discriminative information.

### B. Positives and Negatives Sampling

As the name suggests, contrastive learning requires to obtain two opposing powers: for a given anchor sample, one power is to pull the anchor closer in representation space to other samples, which is known as the positive; while the other power is to push the anchor farther away from other samples, which is known as the negatives. To identify the positive and negative samples, we rely on the constructed global and local memory banks $\mathcal{M}^g$ and $\mathcal{M}^l$ to compute pairwise similarity of samples, and apply two well-designed similarity metrics in the literature to this end.

Specifically, for an anchor image $\mathbf{x}_i$, we learn its global feature $\mathbf{v}_i^g$ and local feature $\{\mathbf{v}_{i,j}^l\}_{j=1}^{N_l}$ in the way as described in last subsection. The similarity between $\mathbf{x}_i$ and another image $\mathbf{x}_j$ are calculated by measuring Euclidean distance of feature vectors and keys of dual dictionary $\mathcal{M}^g$ and $\mathcal{M}^l$. More precisely, we define the global and local distances as $\mathcal{S}^g$ and $\mathcal{S}^l$ respectively as follows:

$$\mathcal{S}^g(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{v}_i^g - \mathcal{M}^g[j]\|_2 \tag{8}$$

and

$$\mathcal{S}^l(\mathbf{x}_i, \mathbf{x}_j) = \frac{\sum_{k=1}^{N_l} \|\mathbf{v}_{i,k}^l - \mathcal{M}^l[j,k]\|_2}{N_l} \tag{9}$$

Moreover, to encourage the consistency of images of the same person captured by different cameras, we add the cross-camera encouragement term (CCE) proposed in [7] as a part of the similarity metric. Set the camera IDs of person images $\mathbf{x}_i$ and $\mathbf{x}_j$ as $c_i$ and $c_j$ respectively, CCE is defined as

$$\mathbb{CCE}(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} \lambda_c & c_i = c_j \\ 0 & c_i \neq c_j \end{cases} \tag{10}$$

Finally, the total similarity metric $\mathcal{S}$ between $\mathbf{x}_i$ and $\mathbf{x}_j$ is formulated as:

$$\mathcal{S}(\mathbf{x}_i, \mathbf{x}_j) = \beta \mathcal{S}^g(\mathbf{x}_i, \mathbf{x}_j) + (1 - \beta)\mathcal{S}^l(\mathbf{x}_i, \mathbf{x}_j) + \mathbb{CCE}(\mathbf{x}_i, \mathbf{x}_j) \tag{11}$$

where $\beta$ is the trade-off parameter that balances the contribution of global and local similarity and is set 0.5.

According to the defined total similarity metric, we perform positive and negative sampling. We rank the similarity order between the anchor and all training samples, according to which we divide the training set into three subsets: *similar set*, *borderline set*, and *dissimilar set*. Considering that each person may have multiple images in dataset, we choose to sample multiple positives for the anchor, as opposed to SimCLR [8] and MoCo [9] that use only a single positive to define the contrastive loss. Specifically, we consider samples in *similar set* as the positives, whose index sets are denoted as $\mathcal{K}_+ \in \mathbb{R}^{N_+}$. Moreover, we propose to select samples that are plausibly similar to the anchor as the negatives, so as to improve the discrimination ability of representation learning. Specifically, we consider in *borderline set* as the negatives, whose index sets are denoted as $\mathcal{K}_- \in \mathbb{R}^{N_-}$. This is different from the conventional contrastive learning strategies that take all samples except the positive as the negatives [8, 9]. And it is also different from an intuitive idea that chooses samples in *dissimilar set* as the negatives. Since samples in *dissimilar set* is already differentiable enough, learning on them cannot improve the discrimination ability of the model significantly. Instead, we enforce the model to consider samples in *borderline set*, which are hard cases with respect to discrimination. We refer to as the proposed manner as selective contrastive learning.

### C. Contrastive Loss and Optimization

With the identified postives and negatives, we define the following contrastive loss function with respect to the mixture

---

**Algorithm 1** Network Training Flow.

---

**Require:** Training set $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N\}$; Global memory bank $\mathcal{M}^g$; Local memory bank $\mathcal{M}^l$; Mixture memory bank $\mathcal{M}^t$; network with parameter $\mathcal{F}(\boldsymbol{\theta}; \mathbf{x}_i)$ ; initial learning rate $\boldsymbol{\gamma}$ ;

**Ensure:** The optimal parameters $\boldsymbol{\theta}^*$

1: **for** $j$ = 1: num_epochs **do**
2:    Pick up training batch set $\{\mathbf{x}_i\}$.
3:    **for** $m$ = 1:n **do**
4:       $\mathbf{f}_i^g = \mathbb{AP}(\mathbf{f}_i)$; $\mathbf{f}_{i,j}^l = \mathbb{AP}(\mathbf{f}_{i,j}), j = 1, \cdots, N_l$
5:       $\mathbf{v}_i^g = \mathbf{P}\left(\mathbb{AP}(\mathbf{f}_i)\right)$
6:       $\mathbf{v}_{i,j}^l = \mathbf{P}\left(\mathbb{AP}(\mathbf{f}_{i,j})\right); j = 1, \cdots, N_l$
7:       $\mathbf{v}_i^l = \mathbb{CONCAT}\left(\{\mathbb{AP}(\mathbf{f}_{i,j})\}_{j=1}^{N_l}\right)$
8:       **if** $i < N_e$ **then**
9:          $\mathcal{L}^g = \mathcal{L}^{init}(\mathbf{v}_i^g | \mathcal{M}^t)$
10:        $\mathcal{L}^l = \mathcal{L}^{init}(\mathbf{v}_i^l | \mathcal{M}^t)$
11:       **else**
12:          Calculate similarities between the anchor image and the other images.
13:          **for** $j$ = 1:N **do**
14:             **if** $j \neq i$ **then**
15:                $\mathcal{S}^g(x_i, x_j) = \|\mathbf{v}_i^g - \mathcal{M}^g[j]\|_2$
16:                $\mathcal{S}^l(x_i, x_j) = \frac{\sum_{k=1}^{N_l}\|\mathbf{v}_{i,k}^l - \mathcal{M}^l[j,k]\|_2}{N_l}$
17:                $\mathcal{S}(\mathbf{x}_i, \mathbf{x}_j) = \beta\mathcal{S}^g(\mathbf{x}_i, \mathbf{x}_j) + (1 - \beta)\mathcal{S}^l(\mathbf{x}_i, \mathbf{x}_j) + \mathbb{CCE}(\mathbf{x}_i, \mathbf{x}_j)$
18:             **end if**
19:          **end for**
20:          Generate set with descending similarity $[\mathbf{x}_{j_1}, \mathbf{x}_{j_2}, ..., \mathbf{x}_{j_{N-1}}]$
21:          Sample positives $[\mathbf{x}_i, \mathbf{x}_{j_1}, \mathbf{x}_{j_2}, ..., \mathbf{x}_{j_{N_+}}]$
22:          Sample negatives $[\mathbf{x}_{j_{N_+}}, ..., \mathbf{x}_{j_{N_- + N_+}}]$
23:          $\mathcal{L}^g = \mathcal{L}(\mathbf{v}_i^g | \mathcal{M}^t)$
24:          $\mathcal{L}^l = \mathcal{L}(\mathbf{v}_i^l | \mathcal{M}^t)$
25:       **end if**
26:       $\mathcal{L}^t(\boldsymbol{\theta}) = (1 - \lambda_p)\mathcal{L}^g + \lambda_p\mathcal{L}^l$
27:    **end for**;
28:    $\boldsymbol{\theta} = \boldsymbol{\theta} - \boldsymbol{\gamma} * \mathbb{SGD}(\nabla_\theta\mathcal{L}^t(\boldsymbol{\theta}))$;
29:    Update $M^g, M^l, M^t$
30: **end for**;
31: $\boldsymbol{\theta}^* = \boldsymbol{\theta}$.

---

memory bank $\mathcal{M}^t$:

$$\mathcal{L}(\mathbf{v}_i | \mathcal{M}^t) = -\log \frac{\sum_{k \in \mathcal{K}_+} \exp(\mathbf{v}_i \cdot \mathcal{M}^t[k]/\tau) * \mu_k}{\sum_{k \in \mathcal{K}_+ \cup \mathcal{K}_-} \exp(\mathbf{v}_i \cdot \mathcal{M}^t[k]/\tau)} \quad (12)$$

where $\cdot$ represents dot product, $\tau$ is a temperature hyper-parameter [22]. In order to emphsize the contribution of the most similar positive sample, we introduce the contribution factor $\mu_k$ for positive distance, which is defined as

$$\mu_k = \begin{cases} \lambda_t & k = i \\ \frac{\alpha(1-\lambda_t)}{|\mathcal{K}_+|} & k \in \mathcal{K}_+ \& k \neq i \end{cases} \quad (13)$$

where $\alpha$ is the expanding coefficient and is set 1.75.

According to the loss form in Eq. (12), we calculate the global loss $\mathcal{L}^g$ and the local loss $\mathcal{L}^l$ with the global feature

$\mathbf{v}_i^g$ defined as Eq. (3) and the local feature $\mathbf{v}_i^l$ defined as Eq. (5), respectively:

$$\mathcal{L}^g = \sum_{i=0}^N \mathcal{L}(\mathbf{v}_i^g | \mathcal{M}^t) \quad (14)$$

and

$$\mathcal{L}^l = \sum_{i=0}^N \mathcal{L}(\mathbf{v}_i^l | \mathcal{M}^t) \quad (15)$$

And finally the total contrastive loss is defined as:

$$\mathcal{L}^t(\boldsymbol{\theta}) = (1 - \lambda_p)\mathcal{L}^g + \lambda_p\mathcal{L}^l \quad (16)$$

where $\lambda_p$ is the trade-off parameter that balances the contributions of global and local losses and is set 0.5. In our framework, the parameters of $\mathbf{E}(\cdot)$ and $\mathbf{P}(\cdot)$ are collectively denoted as $\boldsymbol{\theta}$.

The optimal parameter $\boldsymbol{\theta}^*$ of ReID model $\mathcal{F}(\boldsymbol{\theta}, \cdot)$ can be obtained by:

$$\boldsymbol{\theta}^* = \arg\min_{\boldsymbol{\theta}} \mathcal{L}^t(\boldsymbol{\theta}) \quad (17)$$

This minimization problem can be addressed by stochastic gradient descent (SGD):

$$\boldsymbol{\theta} = \boldsymbol{\theta} - \boldsymbol{\gamma} * \mathbb{SGD}(\nabla_\theta\mathcal{L}^t(\boldsymbol{\theta})) \quad (18)$$

where $\boldsymbol{\gamma}$ is the learning rate, and $\mathbb{SGD}(\nabla_\theta\mathcal{L}^t(\boldsymbol{\theta}))$ represents the updated value based on SGD. The whole network training flow is summarized in Algorithm 1.

It is worth noting that, at the beginning of network training, we set a few epochs to initialize the network and the memory banks. During this process, the positive is the anchor itself and the negatives are randomly chosen. The loss function in initialization stage is defined as follows:

$$\mathcal{L}^{init} = -\sum_{i=0}^N \log \frac{\exp(\mathbf{v}_i \cdot \mathcal{M}^t[i]/\tau)}{\sum_{k \in \{i\} \cup \mathcal{K}_-} \exp(\mathbf{v}_i \cdot \mathcal{M}^t[k]/\tau)} \quad (19)$$

### D. Memory Banks Update

Contrastive learning can be thought of as training an encoder for a dictionary look-up task [9]. In our method, as stated above, there are three dictionaries involved, including the global and local memory banks $\mathcal{M}^g$ and $\mathcal{M}^l$ that are jointly exploited to compute pairwise similarity, and the mixture memory bank $\mathcal{M}^t$ that is used to define constrative loss. To facilitate contrastive unsupervised learning, these three dictionaries should be dynamic, i.e., be updated on-the-fly to provide evolutionary keys during training.

In this work, we propose to use different update strategies for the global and local memory banks and the mixture memory bank, considering that they serve for different purposes. Specifically, for $\mathcal{M}^g$ and $\mathcal{M}^l$, we only update the key corresponding to the anchor $\mathbf{x}_i$ by fusing with the newest global and local feature of $\mathbf{x}_i$ respectively:

$$\mathcal{M}^g[i] = \|\frac{\mathcal{M}^g[i] + \mathbf{v}_i^g}{2}\|_2, \quad (20)$$

and

$$\mathcal{M}^l[i] = \|\frac{\mathcal{M}^l[i] + \mathbf{v}_i^l}{2}\|_2, \quad (21)$$

For the mixture memory bank, we will update the key corresponding to the anchor $\mathbf{x}_i$ and its positives. The update strategy is consistent with the construction strategy, as formulated in Eq. (6) and Eq. (7).

### E. Why Using Three Dynamic Dictionaries?

In our framework, three dynamic dictionaries are used, including the global, local and mixture memory banks. We discuss here about the necessity of using these three dictionaries.

In pairwise similarity computation, we use both the global and local memory banks instead of the mixture one. This is because the keys of the mixture memory bank are the fusion of the global and local features. The fusion operation would remove some useful information. In order to preserve fine information, we thus use both the original global and local features, which construct the global and local memory banks.

In defining the contrastive loss, we use the mixture memory bank instead of the global and local ones. For the task of person ReID, we except the model to generate similar feature representation for the samples with same identity and dissimilar representation otherwise. According to this principle, the mixture memory bank is tailored, in which we update the keys corresponding to the positives of the anchor $\mathbf{x}_i$. This would encourage the anchor and its positives have similar representation and update the keys with global features and local features. It also could provide more discriminative ability to pull the similar samples closer and push dissimilar samples apart away.

## IV. EXPERIMENTS

In this section, we provide extensive experimental results to demonstrate the superior performance of our method.

### A. Dataset

We evaluate our method on five widely used image and video person ReID datasets, including:

- *Market1501* [23], which consists of 32,668 images of 1,501 identities under 6 cameras;
- *DukeMTMC-ReID* [24], which contains 1,812 identities and 36,411 images under 8 cameras;
- *MSMT17* [25], which contains 126,411 person images of 4101 identities under 15 cameras. It is a more challenging dataset due to the effect of substantial variations of scene and lighting.
- *DukeMTMC-VideoReID* [26], which is a video-based ReID dataset containing 2,196 tracklets of 702 identities for training, 2,636 tracklets of other 702 identities for testing;
- *MARS* [27], which is a video-based dataset for person ReID containing 17,503 video tracklets of 1,261 identities.

The first three datasets are image-based, and the last two ones are video-based. We test both image and video based to comprehensively evaluate the performance of our method.

### B. Experimental Setting

We follow the same experimental setting as [7]. All experiments are implemented on PyTorch. The input images are resized to 256*128 and we use random horizontal flip as the data argument strategy. We adopt SGD with momentum as 0.9 to optimize the model. The learning rate is set as 1e-3. The training epoch for image-based dataset is set as 50 and for the video-based dataset is set as 60. The batch size is set as 8. For the video-based dataset, we randomly sample four frames during training, and all frames during testing, in the tracklet. We take the average feature of all frames within a tracklet to be the tracklet feature.

In the proposed framework, there are a few hyper-parameters involved. We set $\lambda_c = 0.005$ in the cross-camera encouragement term, $\beta = 0.5$ in the total similarity metric, $\lambda_t = 0.5$ and $\alpha = 1.75$ in the contribution factor of contrastive loss, the positive samples number $N_+ = 7$ and the negative samples number $N_- = 500$, and $\lambda_p = 0.5$ in the total contrastive loss. The experimental analysis about hyper-parameters setting can be found in Fig. 2, which is conducted on Market-1501.

### C. Comparison with the state-of-the-arts

Our method is comprehensively compared with state-of-the-art unsupervised domain adaption based (UDA), one example based (OneEx) and unsupervised learning based (Unsup) methods. The comparison is conducted on both image-based and video-based datasets.

- **Evaluation on Image-based Datasets:** The comparisons with the state-of-the-art algorithms are conducted on Market-1501, DukeMTMC-ReID and MSMT17, as shown in Table I. It can be found that, under the same setting, our method achieves the best accuracy on both Market-1501 and DukeMTMC-ReID among the 14 compared methods with respect to four performance evaluation metrics: rank-1, rank-5, rank-10 and mAP.
  On Market-1501, we obtain the best performance among the compared methods with rank-1 = 82.2% and mAP = 54.4%. SNR [13], SSLR [7], MMCL [6] and TSSL [14] are four latest methods on unsupervised person ReID. Compared with SNR, we achieve 15.5% and 20.5% improvement on rank-1 accuracy and mAP; compared with SSLR, the gains are 10.5% and 16.6% respectively; compared with MMCL, the gains are 15.6% and 19.1% respectively; compared with TSSL, the gains are 11% and 11.1% respectively. On DukeMTMC-ReID, our method also works the best and achieves accuracy improvement by a large margin. Compared with the second best performed method, we achieve 7.7% and 8.7% improvement on rank-1 accuracy and mAP respectively. On MSMT17, as shown in Table II, compared with the best UDA method SSG, we achieve 9.2% improvement on rank-1. Compared with fully unsupervised method MMCL, we achieve 6.0% and 2.1% improvement on rank-1 and mAP respectively.
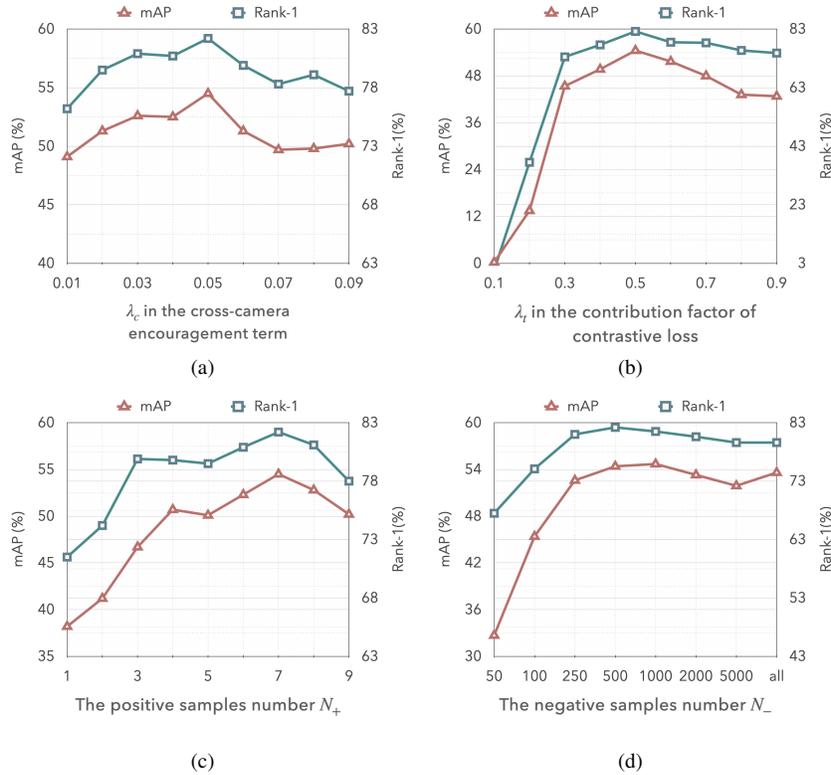
Fig. 2: Experimental analysis about hyper-parameters setting

| Method | Reference | Setting | Market-1501 | | | | | DukeMTMC-reID | | | | |
|--------|-----------|---------|--------|--------|--------|---------|------|--------|--------|--------|---------|------|
| | | | Source | Rank-1 | Rank-5 | Rank-10 | mAP | Source | Rank-1 | Rank-5 | Rank-10 | mAP |
| EUG [16] | CVPR'2018 | OneEx | Market | 49.8 | 66.4 | 72.7 | 22.5 | Duke | 45.2 | 59.2 | 63.4 | 24.5 |
| ATNet [3] | CVPR'2019 | UDA | Duke | 55.7 | 73.2 | 74.9 | 25.6 | Market | 45.1 | 59.5 | 64.2 | 24.9 |
| ProLearn [17] | TIP'2019 | OneEx | Market | 55.8 | 72.3 | 78.4 | 26.2 | Duke | 48.8 | 63.4 | 68.4 | 28.5 |
| SPGAN [2] | CVPR'2018 | UDA | Duke | 58.1 | 76.0 | 82.7 | 26.7 | Market | 46.9 | 62.6 | 68.5 | 26.4 |
| TJ-AIDL [1] | CVPR'2018 | UDA | Duke | 58.2 | - | - | 26.5 | Market | 44.3 | - | - | 23.0 |
| BUC [5] | AAAI'2019 | Unsup | None | 61.0 | 71.6 | 76.4 | 30.6 | None | 40.2 | 52.7 | 57.4 | 21.9 |
| HHL [15] | ECCV'2018 | UDA | Duke | 62.2 | 78.8 | 84.0 | 31.4 | Market | 46.9 | 61.0 | 66.7 | 27.2 |
| DBC [28] | BMVC'2019 | Unsup | None | 69.2 | 83.0 | <u>87.8</u> | 41.3 | None | 51.5 | 64.6 | <u>70.1</u> | 30.0 |
| SNR [13] | CVPR'2020 | UDA | Duke | 66.7 | - | - | 33.9 | Market | 55.1 | - | - | 33.6 |
| SSLR [7] | CVPR'2020 | Unsup | None | <u>71.7</u> | <u>83.8</u> | 87.4 | 37.8 | None | 52.5 | 63.5 | 68.9 | 28.6 |
| MMCL [6] | CVPR'2020 | Unsup | None | 66.6 | - | - | 35.3 | None | 58.0 | - | - | 36.3 |
| TSSL [14] | AAAI'2020 | Unsup | None | 71.2 | - | - | <u>43.3</u> | None | <u>62.2</u> | - | - | <u>38.5</u> |
| Ours | This paper | Unsup | None | **82.2** | **89.9** | **92.6** | **54.4** | None | **69.9** | **79.7** | **82.2** | **47.2** |

TABLE I: The evaluation results with respect to rank-k/mAP on image-based dataset Market-1501 and DukeMTMC. The best and the second ones are highlighted by bold and underline.

The impressive performance demonstrates that the proposed selective contrastive learning framework is able to learn a powerful discrimination model.

- **Evaluation on Video-based Datasets:** We further compare our method with the state-of-the-art algorithms on the two video-based datasets: DukeMTMC-VideoReID and MARS. The comparison results are shown in Table III. On DukeMTMC-VideoReID, we obtain rank-1 = 82.3%, mAP = 78.4%, which are best among all methods. Compared with SSLR, the gains are 5.8% and 9.1% respectively; compared with TTSL, the gains are 8.3% and 13.8% respectively. On MARS, we obtain rank-1 = 66.7%, mAP = 46.8%, which are also the best results.

### D. Ablation Study

In this section, we provide ablation study about the two main contributions of this work: selective contrastive learning and joint usage of global and local features. The experiments are conducted on Market-1501 and DukeMTMC-ReID. The results are reported in Table IV and Table V.

- **Influence of local and global features to the final performance:** We first study the role of global and local features to the final performance. We investigate three scenarios:
  - **Global feature only**: in this case, we use global feature only in pairwise similarity computation and contrastive loss. This is achieved by setting $\beta = 1$

| Method | Reference | Setting | MSMT17 | | | | |
|--------|-----------|---------|--------|--------|--------|---------|-----|
| | | | Source | Rank-1 | Rank-5 | Rank-10 | mAP |
| PTGAN [25] | CVPR'2018 | UDA | Market | 10.2 | - | 24.4 | 2.9 |
| ECN [29] | CVPR'2020 | UDA | Market | 25.3 | 36.3 | 42.1 | 8.5 |
| SSG [30] | ICCV'2019 | UDA | Market | 31.6 | - | 49.6 | 13.2 |
| PTGAN [25] | CVPR'2018 | UDA | Duke | 11.8 | - | 27.4 | 3.3 |
| ECN [29] | CVPR'2020 | UDA | Duke | 30.2 | 41.5 | 46.8 | 10.2 |
| SSG [30] | ICCV'2019 | UDA | Duke | 32.2 | - | <u>51.2</u> | <u>13.3</u> |
| MMCL [31] | CVPR'2020 | Unsup | None | <u>35.4</u> | <u>44.8</u> | 49.8 | 11.2 |
| Ours | This paper | Unsup | None | **41.4** | **53.6** | **58.7** | **13.3** |

TABLE II: The evaluation results with respect to rank-k/mAP on image-based dataset MSMT17. The best and the second ones are highlighted by bold and underline.

| Method | Reference | Setting | DukeMTMC-VideoReID | | | | MARS | | | |
|--------|-----------|---------|--------|--------|---------|-----|--------|--------|---------|-----|
| | | | Rank-1 | Rank-5 | Rank-10 | mAP | Rank-1 | Rank-5 | Rank-10 | mAP |
| RACE [18] | ECCV'2018 | OneEx | - | - | - | - | 43.2 | 57.1 | 62.1 | 24.5 |
| DAL [32] | BMVC'2018 | Unsup | - | - | - | - | 49.3 | 65.9 | 72.2 | 23.0 |
| BUC [5] | AAAI'2019 | Unsup | 76.2 | 88.3 | 91.0 | 68.3 | 57.9 | 72.3 | 75.9 | 34.7 |
| EUG [33] | CVPR'2018 | Unsup | 72.7 | 84.1 | - | 63.2 | 62.6 | 74.9 | - | 42.4 |
| SSLR [7] | CVPR'2020 | Unsup | <u>76.4</u> | <u>88.7</u> | 91.0 | <u>69.3</u> | <u>62.8</u> | **77.2** | **80.1** | 43.6 |
| TTSL [14] | AAAI'2020 | Unsup | 73.9 | - | - | 64.6 | 56.3 | - | - | 30.5 |
| Ours | This paper | Unsup | **82.2** | **93.2** | **95.2** | **78.4** | **66.6** | <u>77.0</u> | <u>79.8</u> | **46.6** |

TABLE III: The evaluation results with respect to rank-k/mAP on video-based dataset DukeMTMC-VideoReID and MARS. The best and the second ones are highlighted by bold and underline.

| Scenarios | Market-1501 | | DukeMTMC | |
|-----------|------|--------|------|--------|
| | mAP | Rank-1 | mAP | Rank-1 |
| Global feature only | 38.3 | 62.0 | 26.9 | 41.2 |
| Local feature only | 42.4 | 70.8 | 34.1 | 57.8 |
| Our joint usage | **54.4** | **82.2** | **47.2** | **69.9** |

TABLE IV: The ablation study about the influence of local and global features to the final performance

| Scenarios | Market-1501 | | DukeMTMC | |
|-----------|------|--------|------|--------|
| | mAP | Rank-1 | mAP | Rank-1 |
| $N_+ = 1, N_- = all$ | 37.1 | 70.9 | 37.9 | 60.3 |
| $N_+ = 7, N_- = all$ | 53.6 | 79.6 | 42.4 | 65.8 |
| Ours ($N_+ = 7, N_- = 500$) | **54.4** | **82.2** | **47.2** | **69.9** |

TABLE V: The ablation study about the influence of positives and negatives to the final performance

and $\lambda_p = 0$ in Eq. (11) and Eq. (16), respectively. The mixture memory bank used in contrastive loss stores global features only.

– **Local feature only**: in this case, we use local feature for these purposes. This is achieved by setting $\beta = 0$ and $\lambda_p = 1$. The mixture memory bank stores local features only.

– **Our joint usage**. This is what we do in this work, *i.e.*, jointly using global and local features.

From Table IV, it can be found that, the joint usage strategy shows the best performance, which achieves improvements with respect to mAP and rank-1 accuracy with a large margin compared with strategies with global or local feature only. This demonstrates that our proposal of jointly using global and local features in dictionary construction is reasonable and works well.

• **Influence of positives and negatives to the final perfor-**

**mance:** In our proposed selective contrastive learning, we leverage multiple positives and selective negatives in contrastive loss definition. Here we study the role of multiple positives and selected negatives to the final performance, and investigate the following three scenarios:

– $N_+ = 1, N_- = all$: This case means that, a single positive is used, and all samples except the positive are used as the negatives, which is what MoCo does [9].

– $N_+ = 7, N_- = all$: This case means that multiple positives are used and all samples except the positives are used as the negatives.

– $N_+ = 7, N_- = 500$: This is what our method does. We use 7 similar samples as positives, and 500 borderline similar samples as negatives.

From Table V, it can be found that, when $N_- = all$, using $N_+ = 7$ positives can significantly improve rank-1 and mAP accuracy compared with using a single positive. This result demonstrates that our proposal of using multiple positives in contrastive learning is useful. Moreover, when $N_+ = 7$, using selected $N_- = 500$ negatives achieves higher rank-1 and mAP accuracy than that using $N_- = all$ negatives. This result demonstrates that our proposed selective choice strategy of negatives is also helpful for contrastive learning.

## V. CONCLUSION

In this work, we presented a novel unsupervised person ReID scheme based on selective contrastive learning. We propose to use multiple positives and adaptively selected negatives for defining the contrastive loss, so as to learn a feature embedding model with stronger discriminative representation ability. We define three dynamic dictionaries for

pairwise similarity computation and contrastive loss definition, which jointly leverage the global and local discriminative information. Experimental results show that our proposed method outperforms the state-of-the-art algorithms.

## REFERENCES

[1] J. Wang, X. Zhu, S. Gong, and W. Li, "Transferable joint attribute-identity deep learning for unsupervised person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2275–2284.

[2] W. Deng, L. Zheng, Q. Ye, G. Kang, Y. Yang, and J. Jiao, "Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 994–1003.

[3] J. Liu, Z.-J. Zha, D. Chen, R. Hong, and M. Wang, "Adaptive transfer network for cross-domain person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 7202–7211.

[4] Y. Ge, F. Zhu, D. Chen, R. Zhao, and H. Li, "Self-paced contrastive learning with hybrid memory for domain adaptive object re-id," in *Advances in Neural Information Processing Systems*, 2020.

[5] Y. Lin, X. Dong, L. Zheng, Y. Yan, and Y. Yang, "A bottom-up clustering approach to unsupervised person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8738–8745.

[6] D. Wang and S. Zhang, "Unsupervised person re-identification via multi-label classification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 10 981–10 990.

[7] Y. Lin, L. Xie, Y. Wu, C. Yan, and Q. Tian, "Unsupervised person re-identification via softened similarity learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3390–3399.

[8] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020.

[9] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9729–9738.

[10] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang, "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 480–496.

[11] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3431–3440.

[12] L. Zhao, X. Li, Y. Zhuang, and J. Wang, "Deeply-learned part-aligned representations for person re-identification,"

[13] X. Jin, C. Lan, W. Zeng, Z. Chen, and L. Zhang, "Style normalization and restitution for generalizable person re-identification," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3140–3149.

[14] G. Wu, X. Zhu, and S. Gong, "Tracklet self-supervised learning for unsupervised person re-identification," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 7, pp. 12 362–12 369, 2020.

[15] Z. Zhong, L. Zheng, S. Li, and Y. Yang, "Generalizing a person retrieval model hetero-and homogeneously," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 172–188.

[16] Y. Wu, Y. Lin, X. Dong, Y. Yan, and Y. Yang, "Exploit the unknown gradually: One-shot video-based person re-identification by stepwise learning," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[17] Y. Wu, Y. Lin, X. Dong, Y. Yan, W. Bian, and Y. Yang, "Progressive learning for person re-identification with one example," *IEEE Transactions on Image Processing*, vol. PP, no. 6, pp. 1–1, 2019.

[18] M. Ye, X. Lan, and P. C. Yuen, "Robust anchor embedding for unsupervised video person re-identification in the wild," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 170–186.

[19] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2197–2206.

[20] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1116–1124.

[21] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *The International Conference on Learning Representations*, 2018.

[22] Z. Wu, Y. Xiong, S. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance-level discrimination," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[23] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *2015 IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 1116–1124.

[24] E. Ristani, F. Solera, R. Zou, R. Cucchiara, and C. Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," in *European Conference on Computer Vision workshop on Benchmarking Multi-Target Tracking*, 2016.

[25] L. Wei, S. Zhang, W. Gao, and Q. Tian, "Person transfer gan to bridge domain gap for person re-identification," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.

in *The IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3219–3228.

[26] Y. Wu, Y. Lin, X. Dong, Y. Yan, W. Ouyang, and Y. Yang, "Exploit the unknown gradually: One-shot video-based person re-identification by stepwise learning," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5177–5186.

[27] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, and Q. Tian, "Mars: A video benchmark for large-scale person re-identification," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 868–884.

[28] G. Ding, S. H. Khan, and Z. Tang, "Dispersion based clustering for unsupervised person re-identification." in *BMVC*, 2019, p. 264.

[29] Z. Zhong, L. Zheng, Z. Luo, S. Li, and Y. Yang, "Invariance matters: Exemplar memory for domain adaptive person re-identification," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[30] Y. Fu, Y. Wei, G. Wang, Y. Zhou, H. Shi, and T. Huang, "Self-similarity grouping: A simple unsupervised cross domain adaptation approach for person re-identification," in *2019 International Conference on Computer Vision (ICCV)*, 2018.

[31] D. Wang and S. Zhang, "Unsupervised person re-identification via multi-label classification," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.

[32] Y. Chen, X. Zhu, and S. Gong, "Deep association learning for unsupervised video person re-identification," *arXiv preprint arXiv:1808.07301*, 2018.

[33] Y. Wu, Y. Lin, X. Dong, Y. Yan, W. Ouyang, and Y. Yang, "Exploit the unknown gradually: One-shot video-based person re-identification by stepwise learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5177–5186.