

# Learn by Guessing: Multi-step Pseudo-label Refinement for Person Re-Identification

Tiago De C. G. Pereira<sup>a</sup> and Teofilo E. De Campos<sup>b</sup>

*Departamento de Ciência da Computação, Universidade de Brasília - UnB, Brasília-DF, Brazil*

**Keywords:** Unsupervised Domain Adaptation, Person Re-Identification, Pseudo-labels, Deep Learning.

**Abstract:** Unsupervised Domain Adaptation (UDA) methods for person Re-Identification (Re-ID) rely on target domain samples to model the marginal distribution of the data. To deal with the lack of target domain labels, UDA methods leverage information from labeled source samples and unlabeled target samples. A promising approach relies on the use of unsupervised learning as part of the pipeline, such as clustering methods. The quality of the clusters clearly plays a major role in methods performance, but this point has been overlooked. In this work, we propose a multi-step pseudo-label refinement method to select the best possible clusters and keep improving them so that these clusters become closer to the class divisions without knowledge of the class labels. Our refinement method includes a cluster selection strategy and a camera-based normalization method which reduces the within-domain variations caused by the use of multiple cameras in person Re-ID. This allows our method to reach state-of-the-art UDA results on DukeMTMC→Market1501 (source→target). We surpass state-of-the-art for UDA Re-ID by 3.4% on Market1501→DukeMTMC datasets, which is a more challenging adaptation setup because the target domain (DukeMTMC) has eight distinct cameras. Furthermore, the camera-based normalization method causes a significant reduction in the number of iterations required for training convergence.

## 1 INTRODUCTION


Person re-identification (Re-ID) aims at matching person images from different non-overlapping cameras views. This is an essential feature for diverse real world challenges, such as smart cities (Zhang and Yu, 2018), intelligent video surveillance (Wang, 2013), suspicious action recognition (Wei Niu et al., 2004) and pedestrian retrieval (Sun et al., 2017).


With all these popular possible applications, there is a clear demand for robust Re-ID systems in the industry. Academic research groups have achieved remarkable in-domain results on popular person Re-ID datasets such as Market1501 (Zheng et al., 2015) and DukeMTMC-reID (Zheng et al., 2017b). Despite these advances, there is still a dichotomy between the success in academic results versus the industrial application. This is because the best academic results e.g. (Wang et al., 2019; Luo et al., 2020; Zhou et al., 2020) are based on supervised methods that require a huge amount of annotated data for their training. Gathering data is not a problem nowadays, as CCTV

systems are omnipresent. However, annotating images is a very expensive and tedious task that requires a lot of manual work.

The use of pre-trained state-of-the-art Re-ID models usually leads to disappointing results because each group of cameras has distinct characteristics, such as illumination, resolution, noise level, orientation, pose, distance, focal length, amount of people's motion as well as factors that influence the appearance of people, such as ethnicity, type of location (e.g. leisure vs work places) and weather conditions.

Some methods have been proposed to reduce this gap and unlock Re-ID systems for real world problems. That is the case of domain invariant models (Jin et al., 2020; Song et al., 2019; Jia et al., 2019), which have the ambitious goal of being applicable to any domain even if no samples are given from some of the potential target domains. Although domain invariance is indeed the key to creating a widely applicable method, such methods usually do not outperform methods in which unlabeled target domain samples are given. Since gathering unlabeled samples is a virtually effortless task, the need for domain invariant methods is not so urgent. A common alternative is the

<sup>a</sup>  <https://orcid.org/0000-0002-9200-9795>

<sup>b</sup>  <https://orcid.org/0000-0001-6172-0229>

use of Generative Adversarial Networks (GANs) to align domains, allowing the model to perform well in a target domain even if supervised training is only performed in the source domain (Zhai et al., 2020; Zhong et al., 2018; Deng et al., 2018; Liu et al., 2019). In addition, there are Unsupervised Domain Adaptation (UDA) methods which have been achieving notable results in cross-domain person Re-ID. These methods typically rely on a process of target domain pseudo-labels generation. This allows them to use actual target domain images without previous annotation (Lin et al., 2020; Ge et al., 2020; Zeng et al., 2020; Zou et al., 2020; Pereira and de Campos, 2020; Fan et al., 2018; Fu et al., 2019).

In this work, we dive deep in the UDA Re-ID setup utilizing pseudo-labels to enhance models performance in target domain. The quality of pseudo-labels clearly is essential for the performance of this kind of method. However, pseudo-labels are expected to be noisy in this scenario. Many methods have use soft cost functions to deal with this noise, however we believe that cleaning and improving pseudo-labels is key to achieve high performance. We therefore focus our work in two main points: camera-based normalization, which we observed to be key to reduce domain variance; and a novel clusters selection strategy. The latter removes outlying clusters and generate pseudo-labels with important characteristics to help model convergence. This strategy aims to generate clusters which are dense and each contain samples of one person captured from the view of multiple cameras.

Enhancing cluster quality has been overlooked by methods based on pseudo-labels and this has certainly held back many methods. To evaluate our proposal we work with the most popular cross-domain dataset in unsupervised Re-ID works: Market1501 and DukeMTMC. Our main contribution is a multi-step pseudo-label refinement that keeps cleaning and improving the predicted target domain label space to enhance model performance without the burden of annotating data. Further to proposing a new pipeline, we introduce strategies to build and select clusters in a way that maximizes the model’s generalization ability and its potential to transfer learning to new Re-ID datasets where the labels are unknown. Our method achieves UDA Re-ID state-of-art for DukeMTMC  $\rightarrow$  Market1501 and significantly pushes state-of-the-art for Market1501  $\rightarrow$  DukeMTMC, improving results in 3.4% w.r.t. the best results we are aware of. We achieve state-of-the-art results without any post-processing methods, however we are aware that re-ranking algorithms are helpful for metric learning tasks. We thus evaluate our model using k-reciprocal

encoding re-ranking (Zhong et al., 2017) and improve our results by further 2.1% and 2.9% for DukeMTMC and Market1501, respectively.

## 2 RELATED WORKS

Person Re-ID has been a trending computer vision research topic. There are two main directions for person Re-ID research: **a)** supervised person Re-ID, that aims at creating the best possible models for in-domain Re-ID and **b)** unsupervised domain adaptation (UDA) Re-ID focusing on the Re-ID task in which a model trained in a source dataset is adapted to another dataset, where the labels are not known. The latter is sometimes referred to as cross-domain Re-ID. In this field, each person Re-ID dataset has images captured from multiple cameras and the dataset as a whole is assumed to be one domain. Although domain adaptation techniques can be applied within a single dataset, i.e., to adapt samples from one camera view to another, we focus on the problem of adapting between different datasets. This setting is more related to a real system deploying setting because training can be done using a dataset containing several viewpoints, but the deployment scenario is a different set of data where it is easy to capture unlabeled samples but labeled samples are not available.

**Generalizable Person Re-ID.** (Jin et al., 2020; Song et al., 2019; Jia et al., 2019; Zhuang et al., 2020) pursue models that are domain invariant, so they can perform well in diverse Re-ID datasets without the need of any adaptation. That is an interesting approach for real world Re-ID challenges, although it still does not perform as well as in-domain person Re-ID. To achieve a domain invariant model, (Zhuang et al., 2020) propose to replace all batch normalization layers of a deep CNN by camera batch normalization layers. These layers learn to apply batch normalization for each camera reducing the within-domain camera variance, this also helps the model to learn camera invariant features that are more robust to domain changes. (Jin et al., 2020) propose a Style Normalization and Restitution (SNR) module that firstly alleviates camera style variations and then restores identity relevant features that may have been discarded in the style normalization step, reducing the influence of camera style change on feature vectors.

**GAN-based Person Re-ID.** (Zhai et al., 2020; Zhong et al., 2018; Deng et al., 2018; Liu et al., 2019; Wei et al., 2018; Zou et al., 2020) have been widely used to reduce the domain gap in Re-ID datasets. (Deng et al., 2018) use cycleGAN to transfer the style from an unlabeled target domain to a labeled source do-

main, so they leverage source domain annotations to apply their trained model to images that are more similar to the source domain ones. (Zhai et al., 2020) used GANs to augment the target domain training data, so they could create images that preserved the person ID and that simulates other camera views at the same time.

**Pseudo-Labels Generation for Person Re-ID.** (Pereira and de Campos, 2020; Ge et al., 2020; Zeng et al., 2020; Fu et al., 2019; Zou et al., 2020; Zhai et al., 2020; Fan et al., 2018; Lin et al., 2020) predict the label space for an unlabeled target domain, assume those predictions are correct and use then to fine-tune a model previously trained on source domain. This approach has shown remarkable results and is the idea behind current state-of-the-art UDA Re-ID methods. The drawback with pseudo-labels is that if the domains are not similar enough, they can lead to negative transfer, because the labeling noise might be too high. To deal with that, (Ge et al., 2020) propose a soft softmax-triplet loss to leverage from pseudo-labels without overfitting their model. (Zeng et al., 2020) propose a hierarchical clustering method to reduce the influence of outliers and use a batch hard triplet loss to bring outliers closer to interesting regions so they could be used later on.

We believe that a generalizable Re-ID model is pre-requisite for a strong UDA method. For this reason we adopt IBN-Net50-a as our backbone and apply a camera guided feature normalization in target domain to reduce the domain gap. We also understand the importance of cleaning the pseudo-labels, which is what motivate us to a clustering algorithm with outlier detection and to propose a clustering selection step to feed our model only with data that is predicted to be more reliable. Next section details our methods.

### 3 METHODOLOGY

In this Section we present and discuss all steps that compose our method in deep details. In §3.1 we discuss the commonly used backbones for Re-ID and in §3.2 we present the concept of progressive learning. Then, we review some clustering techniques in §3.3 and how to generate robust clusters in §3.4 and §3.5. Finally, in §3.6 we explain how to effectively combine all techniques in our training protocol.

#### 3.1 Backbone

When working in cross domain tasks, the model generalization ability is key to success. Normalization

techniques have a very important role for that.

Nowadays, the typical Re-ID system relies on ResNet (He et al., 2016) as their backbone (usually the ResNet-50 model), which is a safe choice, because Re-ID is a task that requires multiple semantic levels to produce robust embeddings and the residual blocks help to propagate these multiple semantic levels to deeper layers. Also, ResNet is a well studied CNN that lead to a step change in the performance on the ImageNet dataset.

However, the vanilla ResNet has its generalization compromised because it does not include instance-batch normalization. (Pan et al., 2018) proposed IBN-Net50, which replaces batch normalization (BN) layers with instance batch normalization (IBN) layers. The IBN-Net carefully integrates IN and BN as building blocks, significantly increasing its generalization ability. For this reason, we choose it as our backbone. More specifically, we use IBN-Net50-a, which offers a good compromise between model size and performance.

#### 3.2 Progressive Learning

Progressive learning is an iterative technique proposed by (Fan et al., 2018) composed of three parts: **a)** generating target domain pseudo-labels to train the model without labeled data, **b)** fine-tuning the model with the previous generated pseudo-labels and **c)** evaluating the model. This set of steps is iterated until convergence. This approach relates to some classical UDN and Transductive Transfer Learning techniques, such as (FarajiDavar et al., 2017) and (Long et al., 2013) for standard classification tasks.

To get full advantage of progressive learning it is important to generate new pseudo-labels at each iteration, so the model will have new stimulus to keep learning. Therefore it is important that in each step the new pseudo-labels get closer of the real labels. However, if the initial model is not good enough, this leads to negative transfer (Pan and Yang, 2010) and the performance of the system actually degrades as it iterates. However, since target labels are unknown, it is not possible to predict negative transfer.

For this reason, we argue that progressive learning must be coupled with other techniques, such as the method we describe in the next sections, particularly in §3.4 and §3.5. In those sections, we propose to evaluate the reliability of samples and their pseudo-labels based on the confidence of the model. If only reliable samples and their pseudo-labels are used, the model should progressively improve and generate more robust pseudo-labels in the consecutive iterations.

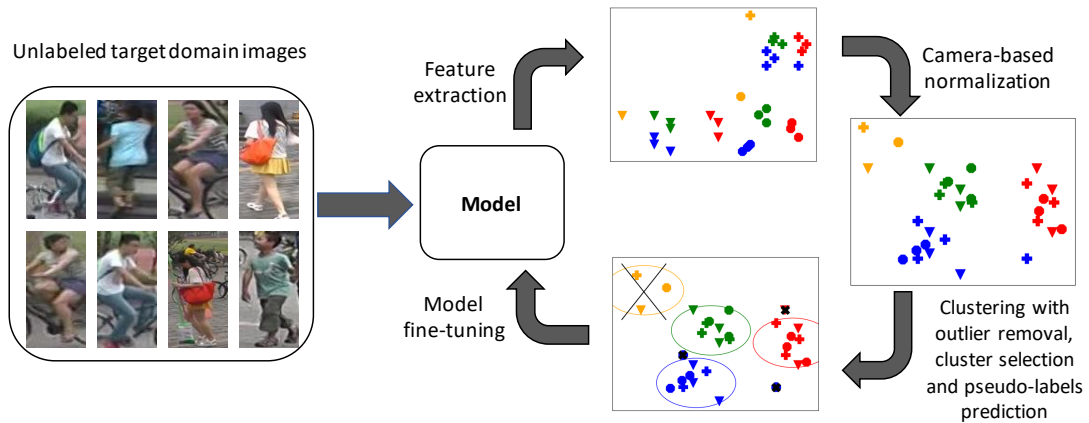


Figure 1: Our Multi-Step Pseudo-Label Refinement pipeline. The proposed method consists of four components: extraction of features from unlabeled target domain images, camera-based normalization, prediction of pseudo-labels with a density-based clustering algorithm, selection of reliable clusters and fine-tuning of the model. The pipeline is cyclical, because at each step it predicts more robust pseudo-labels that offer new information for the model. In the feature space panels, each shape (e.g. triangle, plus signal and circle) represent a camera view and each color represent a person ID.

### 3.3 Clustering Techniques

For standard classification tasks, pseudo-labels generation is direct: it is assumed that the predictions obtained are usually correct and these predictions on the target set are used as pseudo class labels. However, due to the lack of control on the number of classes, person Re-ID is usually approached as a metric learning task. The model prediction is therefore not a label, but a feature vector in a space where samples of the same person are expected to lie closer to each other (and further to samples of different people). Therefore, it is necessary to use clustering algorithms and define each cluster as a pseudo-label (or pseudo person ID).

Therefore, given a target domain  $\mathcal{D}^t$  with  $N$  images  $\{\mathbf{x}_i\}_{i=1}^N$  we need to predict their labels  $\{y_i\}_{i=1}^N$ . Then, we use a model pre-trained on source domain  $\mathcal{D}^s$  to extract the features for each image  $\{\mathbf{x}_i\}_{i=1}^N$  from  $\mathcal{D}^t$  and use a clustering algorithm to group/predict each image label<sup>1</sup>.

#### 3.3.1 K-means

As a first choice, we used the k-means algorithm to cluster our data. The only parameter k-means need is the number of clusters  $k$ . For our experiments, we choose  $k$  using this heuristic:

$$k = \left\lfloor \frac{N}{15} \right\rfloor, \quad (1)$$

where  $N$  is the total number of training images in target domain  $\mathcal{D}^t$ . If all clusters have a balanced number

<sup>1</sup>We use bold math symbols for vectors.

of features (images) this would mean that we are assuming that each person ID in the target domain contains about 15 samples.

There are two problems with k-means for Re-ID: **a)** how to define  $k$  without information about  $\mathcal{D}^t$  and **b)** as stated by (Zeng et al., 2020) k-means does not have an outlier detector, so the outliers may drag the centroids away from denser regions, causing the decision boundaries to shift, potentially cutting through sets of samples that actually belong to the same people.

#### 3.3.2 DBSCAN

As discussed above, k-means is not recommended to generate robust pseudo-labels for UDA Re-ID methods. Therefore, we propose the usage of DBSCAN which is a density-based clustering algorithm designed to deal with large and noisy databases.

In a Domain Adaptation Re-ID scenario we can say that the hard samples are actually noise, so a clustering algorithm that identify them as outliers is fundamental to improve results. Furthermore, when applying Progressive Learning we can leave hard samples out for some iterations and bring them to the pseudo-labeled dataset in later iterations where our model is stronger and the level of confidence in those hard samples is higher.

One important point is that DBSCAN does not require a pre-defined number of clusters (as in k-means), but it requires two parameters: the maximum distance between two samples to determine them as neighbors ( $\epsilon$ ) and the minimum number of samples to consider a region as dense ( $min\_s$ ).



In our experiments, we set  $min\_s = 4$ . As for the parameter  $\epsilon$ , its value depends on the spread of the data. We performed a grid search in an early training step to determine a value that would balance the number of clusters selected and the number of outliers. This lead to  $\epsilon = 0.35$  when DukeMTMC is the target domain and  $\epsilon = 0.42$  when Market1501 is the target domain.

### 3.4 Cluster Selection

Re-ID datasets have disjoint label spaces, that is given a source domain  $\mathcal{D}^s$  and a target domain  $\mathcal{D}^t$  their labels space do not share the same classes, i.e.

$$\{y_i\}^s \neq \{y_j\}^t \forall i,j. \quad (2)$$

As mentioned before, the usual way to deal with this is by approaching Re-ID as a metric-learning task.

Therefore, Re-ID methods typically use triplet loss with batch hard (Hermans et al., 2017) and PK sampling. The PK sampling method consist in selecting  $P$  identities with  $K$  samples from each identity to form a mini-batch in training stage, which leads to the following:

$$batch\_size = P \times K. \quad (3)$$

In this work we used the triplet loss and PK sampling to train our models, so we expect that every person ID has at least  $K$  images. This clustering step therefore ignores clusters with less than  $K$  images.

An important factor for Re-ID models is to learn features that are robust to camera view variations. For that we guarantee that, in the training stage, our model is fed with samples of the same person ID in different cameras. Therefore, we also prune clusters that had images from only one camera view.

### 3.5 Camera-guided Feature Normalization

The high variance present in Re-ID datasets is mainly caused by the different camera views, as each view has its own characteristics. This is why a model trained in a source dataset presents poor results when evaluated in a target dataset (or domain). Normally, Re-ID models learn robust features for known views, but lack the ability to generalize for new unseen views.

(Pereira and de Campos, 2020) realize that this lack of generalization power has a negative impact in pseudo-labels generation. They point that the main reason for that is the fact that, in new unseen cameras, the model tends to cluster images by cameras rather

than clustering images from the same person in different views. The majority of clusters would therefore be ignored in the Cluster Selection step.

(Zhuang et al., 2020) replaced all batch normalization layers by camera batch normalization layers. Although this helped them to reduce the data variance between camera views, they normalize the data only on the source domain. We propose to run this camera feature normalization step before the pseudo-labels step on the target domain training set. By generating pseudo-labels that are normalized by camera information, our method guides the model to learn robust features in the target domain space without the need of changing the model architecture or having additional cost functions.

Camera-guided feature normalization therefore aims to reduce the target domain variance, enhance the model capacity in the target domain and create better pseudo-labels that further will result in a more robust model.

To apply camera guided feature normalization, we first divide all target domain training images  $\{\mathbf{x}_i\}^t$  in  $c$  groups where  $c$  is the number of cameras views in the dataset. Then we extract their features  $\mathbf{f}_i^{(c)}$  with our model and calculate, for each camera  $c$ , its mean  $\mu^{(c)}$  and its standard deviation  $\sigma^{(c)}$ . Finally, each feature is normalized by

$$\bar{\mathbf{f}}_i^{(c)} = \frac{\mathbf{f}_i^{(c)} - \mu^{(c)}}{\sigma^{(c)}}. \quad (4)$$

The normalized features  $\bar{\mathbf{f}}_i^{(c)}$  are then used to generate the pseudo-labels.

## 3.6 Our Training Protocol

### 3.6.1 Baseline

First of all, we train our model in the source domain  $\mathcal{D}^s$  as a baseline. All our models use the IBN-Net50-a as backbone and outputs feature vectors  $\mathbf{f}$  with 2048 dimensions and a logit prediction vectors  $\mathbf{p}$ .

Our loss function has three components: **a**) a batch hard triplet loss ( $\mathcal{L}_{tri}$ ) that maps  $\mathbf{f}$  in an Euclidean vector space, **b**) a center loss ( $\mathcal{L}_c$ ) (Wen et al., 2016) to guarantee cluster compactness and **c**) and a cross entropy label smooth loss ( $\mathcal{L}_{ID}$ ) (Zheng et al., 2017a) that uses the logit vectors  $\mathbf{p}$  to learn a person ID classifier. The smoothed person ID component has been proved to help Re-ID systems even though the training IDs are disjoint from the testing IDs. Furthermore, its soft labels has shown interesting features for UDA Re-ID (Ge et al., 2020). Our loss function is thus given by Equation 5.

$$\mathcal{L} = \mathcal{L}_{tri} + \mathcal{L}_{ID} + 0.005\mathcal{L}_c \quad (5)$$

The weight given to the cross entropy loss is the same that was used in (Luo et al., 2020).

We start our training with pre-trained weights from ImageNet and use the Adam optimizer for 90 epochs with the warm-up learning rate scheduler proposed by (Luo et al., 2020).

### 3.6.2 Unsupervised Domain Adaptation

For unsupervised domain adaptation, we start with the model pre-trained in  $\mathcal{D}^s$  and use it to extract all the features  $f$  from  $\mathcal{D}^t$  training images. Once we have all these features extracted, we separate them by camera and use Equation 4 to normalize them. Then, we use DBSCAN to create general clusters in  $\mathcal{D}^t$  and finally apply our cluster selection strategy to keep only the clusters which are potentially the most reliable ones.

From the selected clusters we create our pseudo-labeled dataset and use it to fine-tune our previous model. Since the domains are different datasets, the person IDs on the pseudo-labeled dataset are always different from those of the source dataset. Additionally, as our progressive learning strategy iterates, pseudo-labels are expected to change. Therefore, it is expected that the cross-entropy loss  $\mathcal{L}_{ID}$  spikes in first iterations, which can destabilize the training process and lead to catastrophic forgetting. To prevent that, we follow the transfer learning strategy of freezing the body of our model for 20 epochs and let the last fully connected layer learn a good enough  $\mathbf{p}$ . Then, we unfreeze our model and complete the fine-tuning following the procedure described in 3.6.1.

After the fine-tuning we evaluate our model on  $\mathcal{D}^t$  and iterate the whole process, according to the progressive learning strategy.

## 4 EXPERIMENTS

### 4.1 Datasets

We performed our experiments on the Market1501 and DukeMTMC datasets, interchanging then as source and target domains to analyze our domain adaptation method. Following the standard in this field, we used cumulative matches curve (CMC) and mean average precision (mAP) as evaluation metrics. **Market1501:** is an outdoors dataset containing images across 6 cameras views where each person appears in at least 2 different cameras. In total there are 32668 images being 12936 images from 751 identities for training and 19732 images from 750 identities for testing.

**DukeMTMC:** was also built using outdoor cameras. It contains 36411 images from 1404 identities. Those images were split as 16522 for training, 2228 for query and 17661 for test gallery.

### 4.2 Comparison with Supervised Learning and Direct Transfer

In Table 1 we compare our baseline results to the direct transfer and to our proposed method. The supervised learning was done using samples and labels from the target domain training samples. Since samples and labels are from the same domain as the test set, this is expected to give results that are better than those of domain adaptation settings. The aim of the supervised learning experiments is to understand the capacity of the model in each dataset.

The Direct transfer method is used to evaluate the domain shift and the model generalization power. It is expected that this setting gives results that are worse than the domain adaptation setting, because no knowledge of the target set is used in the training process. Our method does not focus on being generalizable, we aim to use the source domain knowledge as start and enhance the model’s performance in target domain without any labels. We found it important to present direct transfer results in order to show how much our method enhances over it.

As one can see, our method reaches remarkable results for DukeMTMC as a target dataset. It can be surprising to see that we even surpass the supervised result in 0.3% and 0.5% for CMC rank-1 and mAP, respectively. DukeMTMC is a dataset with a high intra-variance caused by its eight distinct camera views. We believe that the camera-guided normalization applied before the clustering step provided pseudo-labels that were more robust to camera view variations. Therefore, the method was able to learn camera invariant features. It is also likely that by transferring from one dataset to another, our method was less prone to over-fitting than the supervised learning setting.

For Market1501 as a target, our method performed equally well enhancing the direct transfer result in 30.2% and 44.6% for CMC rank-1 and mAP, respectively. However, with lower intra-variance in Market1501 the supervised result is already saturated. Therefore, even though labels from the target set were not used, our methods gives results which are not far below those of the supervised setting.

Table 1: Comparison of our results with results using supervised learning on the target domain (which is expected to give the best results) and direct transfer results, i.e. the use of a model trained on source directly applied to the target domain, without domain adaptation (which is expected to be a lower bound).

Methods	Market1501 → DukeMTMC				DukeMTMC → Market1501			
	Rank-1	Rank-5	Rank-10	mAP	Rank-1	Rank-5	Rank-10	mAP
Supervised	82.4	92.0	94.5	68.8	91.2	92.0	98.4	79.2
Direct Transfer	44.7	60.7	66.4	27.3	58.9	74.3	80.1	29.0
<b>Ours</b>	82.7	90.5	93.5	69.3	89.1	95.8	97.2	73.6

Table 2: Comparison of our results with state-of-art methods in UDA. We highlighted in bold, underline and italic the first, second and third best results, respectively. RR stands for Re-Ranking.

Methods	Market1501 → DukeMTMC				DukeMTMC → Market1501			
	Rank-1	Rank-5	Rank-10	mAP	Rank-1	Rank-5	Rank-10	mAP
ECN (Zhong et al., 2019)	63.3	75.8	80.4	40.4	75.1	87.6	91.6	43.0
CBN (Zhuang et al., 2020) + ECN	68.0	80.0	83.9	44.9	81.7	91.9	94.7	52.0
Theory (Song et al., 2020)	68.4	80.1	83.5	49.0	75.8	89.5	93.2	53.7
PCB-PAST (Zhang et al., 2019)	72.4	-	-	54.3	78.4	-	-	54.6
AD Cluster (Zhai et al., 2020)	72.6	82.5	85.5	54.1	86.7	94.4	96.5	68.3
SSG (Fu et al., 2019)	76.0	85.8	89.3	60.3	86.2	94.6	96.5	68.7
DG-Net++ (Zou et al., 2020)	78.9	87.8	90.4	63.8	82.1	90.2	92.7	61.7
MMT (Ge et al., 2020)	79.3	<i>89.1</i>	<i>92.4</i>	<i>65.7</i>	<i>90.9</i>	<b>96.4</b>	<b>97.9</b>	<i>76.5</i>
<b>Ours</b>	<b>82.7</b>	<b>90.5</b>	<b>93.5</b>	<b>69.3</b>	<i>89.1</i>	<i>95.8</i>	<i>97.2</i>	<i>73.6</i>
<b>Ours + RR (Zhong et al., 2017)</b>	<b>84.8</b>	<b>90.8</b>	<b>93.2</b>	<b>81.2</b>	<b>92.0</b>	<i>95.3</i>	<i>96.6</i>	<b>88.1</b>

### 4.3 Comparison with State-of-art UDA Results

In Table 2 we compare our multi-step pseudo-label refinement method with multiple state-of-the-art Re-ID UDA methods. As one can see, we beat all other methods in DukeMTMC target dataset and push the state-of-the-art by 3.4% and 3.6% for CMC rank-1 and mAP, respectively. For Market1501 we are able to reach second place with a noticeable gap to the third place with an improvement of 2.4% and 4.9% for CMC rank-1 and mAP, respectively.

In addition, our framework have a lightweight architecture when compared to other frameworks that achieve state-of-the-art. MMT (Ge et al., 2020) uses two CNNs so that one generates soft labels for the other. DG-Net++ (Zou et al., 2020) uses a extremely complex framework with GANs and multiple encoders and decoders.

As we approach Re-ID as a metric learning task, re-ranking algorithms have a great impact in the results. Then, we evaluated our model using k-reciprocal encoding re-ranking (Zhong et al., 2017) which combines the original distance with the Jaccard distance in an unsupervised manner. The importance to use a ranking system is shown with an CMC Rank-1 improvement of 2.1% for DukeMTMC and 2.9% in Market1501 when compared to our raw method. Also, re-ranking significantly pushes the mAP performance in 11.9% for DukeMTMC and 14.5% for Market1501.

### 4.4 Ablation Studies

Table 3: The contribution of each method in the model performance evaluated on Market1501 and DukeMTMC-reID datasets. PL means Progressive Learning, CN stands for Camera Guided Normalization and DA for Domain Adaptation. The baseline Resnet-50 results are from (Luo et al., 2020).

Methods	M → D		D → M	
	Rank-1	mAP	Rank-1	mAP
ResNet 50	41.4	25.7	54.3	25.5
+ IBN-Net50-a	44.7	27.3	58.9	29.0
+ DA	52.2	37.1	60.1	34.8
+ PL	52.2	37.1	61.4	35.5
+ Cluster Selection	77.2	61.8	86.5	66.0
+ CN	82.7	69.3	89.1	73.6

Table 3 shows how each technique contributes to our final method performance.

**IBN-Net50-a:** the difference between the original Resnet-50 and the IBN-Net50-a is that the IBN-Net50-a modifies all batch normalization layers so they also take advantage of an instance normalization. This modification enhances the model normalization capacity and generalization power, leading to an improvement on the CMC rank-1 performance improvement of 3.3% in DukeMTMC and 4.6% in Market1501.

**Domain Adaptation:** in Table 3 we call as domain adaptation the use of pseudo-labels from target domain for training. This clustering guided domain adaptation allows our model to train using actual images from target domain, which facilitates the model to learn various aspects of the domain, such as illu-

mination, camera angles, person pose. Learning the characteristics from target domain is a major factor for domain adaptation which becomes evident by the CMC rank-1 improvement of 7.5% in DukeMTMC and 1.2% in Market1501.

**Progressive Learning:** this technique has a great potential to keep improving the model’s performance with new pseudo-labels. However, as we said in Section §3.2 to get full advantage of this technique ones need to guarantee that the pseudo-labels are close to the class divisions. Therefore, this step only gives a significant improvement if associated with the proposed cluster selection technique. In Table 3 results, the progressive learning results were obtained using the raw clusters defined by the clustering algorithm. Then, the model used all the available information in target domain and overfitted to these pseudo-labels. In the next step these clusters tend to be the same and the model does not have a stimulus to learn better features. This is why the progressive learning results on their own do not seem to help.

**Cluster Selection:** this method relies on a continuous improvement on the pseudo-labels. for that, we remove clusters that are unlikely to help improve the model, such as small clusters with less than 4 images and clusters that had images from only one camera view. Using this strategy we can get full advantage of progressive learning and push the model to learn camera view invariant features, since all our pseudo-labels have samples from multiple camera views.

The real contribution of the progressive learning technique is shown alongside the contribution of the cluster selection strategy, because they are complementary techniques. This is certainly the most relevant element of our pipeline, as it leads to a step change in our performance, enhancing the rank-1 CMC performance in 25.0% for DukeMTMC and 25.1% for Market1501.

**Camera Guided Normalization:** learning camera invariant features is essential for person Re-ID, because the person appearance may vary for different cameras and the model has to deal with all types of variations. Since target labels are unknown, when the model extracts features from the target domain, instead of grouping images by the person that appears in them, the feature vectors tend to cluster camera viewpoints. The camera guided normalization helps to reduce this camera shift and align the features from different cameras. This camera alignment allows the cluster method to create better clusters with samples from different cameras. Our cluster selection method thus selects more clusters to be part of the pseudo-label dataset. With this richer and camera invariant pseudo-label dataset, our model has better samples

to learn from and its mAP is improved by 7.5% for DukeMTMC and 7.6% for Market1501.

**Training Efficiency:** the better pseudo-labels which are obtained when applying camera guided normalization speeds up the model convergence. Figure 2 shows how many progressive learning steps were needed to reach convergence with or without camera guided normalization.

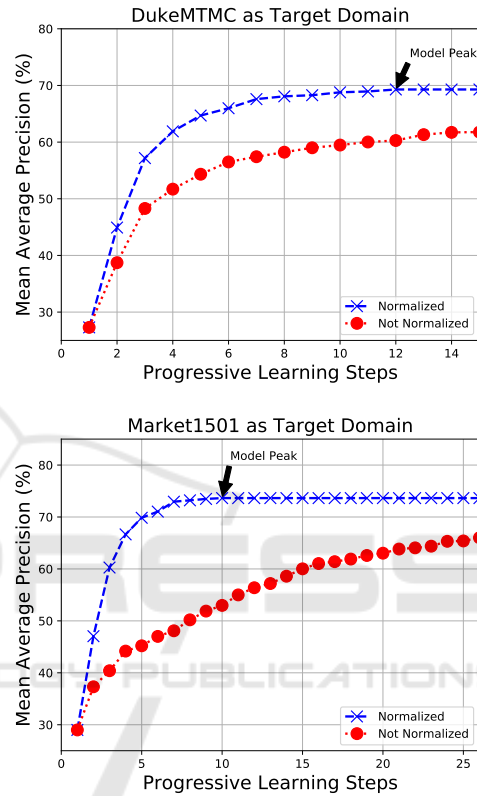


Figure 2: Comparison of progressive learning steps that take for model convergence when using or not camera guided normalization. The black arrow indicates when the model using normalization reached convergence, the points after the black arrow repeat the best model result. The “not normalized” curves reach their peak at the right end of these plots.

Table 4: Comparison between DBSCAN and k-means as the clustering algorithm. After cluster selection, different amounts of samples were removed for each clustering method. This portion (in %) is shown in the last columns.

Method	M → D				Portion
	Rank-1	Rank-5	Rank-10	mAP	
k-means	77.7	87.5	90.8	63.1	85.5
DBSCAN	82.7	90.5	93.5	69.3	69.7

Method	D → M				Portion
	Rank-1	Rank-5	Rank-10	mAP	
k-means	87.0	94.7	96.9	65.9	95.9
DBSCAN	89.1	95.8	97.2	73.6	79.9



**Clustering Methods:** we ran our multi-step pseudo-label refinement method with two different clustering algorithms in its pipeline: k-means and DBSCAN. Table 4 presents the results achieved using each of them and also the portion of training data that was selected for use as pseudo-labels after the cluster selection phases. DBSCAN does not need a fixed number of clusters and has an built-in outlier detector, so it can deal with hard samples better than k-means. For k-means, all samples count, then the hard samples have a negative impact in the quality of the pseudo-labels. The results in Table 4 confirms our hypothesis that it is better to use fewer and less noisy samples.

## 5 CONCLUSIONS

In this work we propose a multi-step pseudo-label refinement method to improve results on Unsupervised Domain Adaptation for Person Re-Identification. We focus on tackling the problem of having noisy pseudo-labels in this task and proposed a pipeline that reduces the shift caused by camera changes as well as techniques for outlier removal and cluster selection. Our method includes DBSCAN clustering algorithm that was designed to perform well in large and noisy databases; a camera-guided normalization step to align features from multiple camera views and allow them to be part of the same clusters; and a smart cluster selection method that creates optimal pseudo-labels for our training setup and keep improving the pseudo-labels at each progressive learning step.

Our method generates a strong label space for target domain without any supervision. We reach state-of-the-art performance on Market1501 as a target dataset and push the state-of-the-art on the challenging DukeMTMC target dataset by 5.5% (or 3.4% without re-ranking). Our work highlights the importance of pseudo-labels refinement with strong normalization techniques. It also takes advantage of a metric learning process and re-ranking (Zhou et al., 2020; Zhong et al., 2017). This combination has clearly proven successful.

One possibility for future work is to investigate the use of re-ranking as part of the clustering step.

## ACKNOWLEDGEMENTS

We acknowledge the support of Fundação de Apoio à Pesquisa do Distrito Federal (FAPDF) through Convênio 07/2019 - project KnEDLe - with support of Finatec (project 6429).

Dr. de Campos also acknowledges the support of CNPq fellowship PQ 314154/2018-3. He currently works at Vicon MotionSystems, Oxford Metrics Group, UK. Mr. Pereira also thanks CyberLabs for supporting his research.

## REFERENCES

- Deng, W., Zheng, L., Ye, Q., Kang, G., Yang, Y., and Jiao, J. (2018). Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Fan, H., Zheng, L., Yan, C., and Yang, Y. (2018). Unsupervised person re-identification: Clustering and fine-tuning. *ACM Transactions on Multimedia Computing, Communications, and Applications TOMM*, 14(4):83:1–83:18.
- FarajiDavar, N., de Campos, T., and Kittler, J. (2017). Adaptive transductive transfer machines: A pipeline for unsupervised domain adaptation. In *Domain Adaptation in Computer Vision Applications, Advances in Computer Vision and Pattern Recognition*, pages 115–132. Springer International.
- Fu, Y., Wei, Y., Wang, G., Zhou, Y., Shi, H., and Huang, T. S. (2019). Self-similarity grouping: A simple unsupervised cross domain adaptation approach for person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Ge, Y., Chen, D., and Li, H. (2020). Mutual mean-teaching: Pseudo label refinery for unsupervised domain adaptation on person re-identification. In *International Conference on Learning Representations*.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hermans, A., Beyer, L., and Leibe, B. (2017). In defense of the triplet loss for person re-identification. Technical Report arXiv:1703.07737, Cornell University Library. <http://arxiv.org/abs/1703.07737>.
- Jia, J., Ruan, Q., and Hospedales, T. M. (2019). Frustratingly easy person re-identification: Generalizing person re-id in practice. In *30th British Machine Vision Conference 2019, BMVC 2019, Cardiff, UK, September 9-12, 2019*, page 117. BMVA Press.
- Jin, X., Lan, C., Zeng, W., Chen, Z., and Zhang, L. (2020). Style normalization and restitution for generalizable person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Lin, Y., Xie, L., Wu, Y., Yan, C., and Tian, Q. (2020). Unsupervised person re-identification via softened similarity learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Liu, J., Zha, Z.-J., Chen, D., Hong, R., and Wang, M. (2019). Adaptive transfer network for cross-

- domain person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Long, M., Wang, J., Ding, G., Sun, J., and Yu, P. S. (2013). Transfer feature learning with joint distribution adaptation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Luo, H., Jiang, W., Gu, Y., Liu, F., Liao, X., Lai, S., and Gu, J. (2020). A strong baseline and batch normalization neck for deep person re-identification. *IEEE Transactions on Multimedia*, 22(10):2597–2609.
- Pan, S. J. and Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359.
- Pan, X., Luo, P., Shi, J., and Tang, X. (2018). Two at once: Enhancing learning and generalization capacities via ibn-net. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- Pereira, T. and de Campos, T. E. (2020). Domain adaptation for person re-identification on new unlabeled data. In *15<sup>th</sup> International Conference on Computer Vision Theory and Applications (VISAPP) - part of VISIGRAPP*, volume 4: VISAPP, pages 695–703.
- Song, J., Yang, Y., Song, Y.-Z., Xiang, T., and Hospedales, T. M. (2019). Generalizable person re-identification by domain-invariant mapping network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Song, L., Wang, C., Zhang, L., Du, B., Zhang, Q., Huang, C., and Wang, X. (2020). Unsupervised domain adaptive re-identification: Theory and practice. *Pattern Recognition*, 102:107173.
- Sun, Y., Zheng, L., Deng, W., and Wang, S. (2017). Svdnet for pedestrian retrieval. In *IEEE International Conference on Computer Vision (ICCV)*, pages 3820–3828.
- Wang, G., Lai, J., Huang, P., and Xie, X. (2019). Spatial-temporal person re-identification. pages 8933–8940.
- Wang, X. (2013). Intelligent multi-camera video surveillance: A review. *Pattern Recognition Letters*, 34(1):3 – 19. Extracting Semantics from Multi-Spectrum Video.
- Wei, L., Zhang, S., Gao, W., and Tian, Q. (2018). Person transfer gan to bridge domain gap for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Wei Niu, Jiao Long, Dan Han, and Yuan-Fang Wang (2004). Human activity detection and recognition for video surveillance. In *IEEE International Conference on Multimedia and Expo (ICME) (IEEE Cat. No.04TH8763)*, volume 1, pages 719–722 Vol.1.
- Wen, Y., Zhang, K., Li, Z., and Qiao, Y. (2016). A discriminative feature learning approach for deep face recognition. In *ECCV (7)*, pages 499–515.
- Zeng, K., Ning, M., Wang, Y., and Guo, Y. (2020). Hierarchical clustering with hard-batch triplet loss for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhai, Y., Lu, S., Ye, Q., Shan, X., Chen, J., Ji, R., and Tian, Y. (2020). Ad-cluster: Augmented discriminative clustering for domain adaptive person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhang, S. and Yu, H. (2018). Person re-identification by multi-camera networks for internet of things in smart cities. *IEEE Access*, 6:76111–76117.
- Zhang, X., Cao, J., Shen, C., and You, M. (2019). Self-training with progressive augmentation for unsupervised cross-domain person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*.
- Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., and Tian, Q. (2015). Scalable person re-identification: A benchmark. In *IEEE International Conference on Computer Vision*.
- Zheng, Z., Zheng, L., and Yang, Y. (2017a). A discriminatively learned cnn embedding for person re-identification. *ACM Trans. Multimedia Comput. Commun. Appl.*, 14(1).
- Zheng, Z., Zheng, L., and Yang, Y. (2017b). Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- Zhong, Z., Zheng, L., Cao, D., and Li, S. (2017). Re-ranking person re-identification with k-reciprocal encoding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhong, Z., Zheng, L., Luo, Z., Li, S., and Yang, Y. (2019). Invariance matters: Exemplar memory for domain adaptive person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhong, Z., Zheng, L., Zheng, Z., Li, S., and Yang, Y. (2018). Camera style adaptation for person re-identification. In *Proc of the IEEE Conf on Computer Vision and Pattern Recognition (CVPR)*.
- Zhou, J., Su, B., and Wu, Y. (2020). Online joint multi-metric adaptation from frequent sharing-subset mining for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhuang, Z., Wei, L., Xie, L., Zhang, T., Zhang, H., Wu, H., Ai, H., and Tian, Q. (2020). Rethinking the distribution gap of person re-identification with camera-based batch normalization. In *ECCV*.
- Zou, Y., Yang, X., Yu, Z., Kumar, B. V. K. V., and Kautz, J. (2020). Joint disentangling and adaptation for cross-domain person re-identification. In *ECCV*.