

Confidence-guided Centroids for Unsupervised Person Re-Identification

Yunqi Miao
University of Warwick

Jiankang Deng
Huawei

Guiguang Ding
Tsinghua University

Jungong Han
Aberystwyth University

Abstract

Unsupervised person re-identification (ReID) aims to train a feature extractor for identity retrieval without exploiting identity labels. Due to the blind trust in imperfect clustering results, the learning is inevitably misled by unreliable pseudo labels. Albeit the pseudo label refinement has been investigated by previous works, they generally leverage auxiliary information such as camera IDs and body part predictions. This work explores the internal characteristics of clusters to refine pseudo labels. To this end, Confidence-Guided Centroids (CGC) are proposed to provide reliable cluster-wise prototypes for feature learning. Since samples with high confidence are exclusively involved in the formation of centroids, the identity information of low-confidence samples, i.e., boundary samples, are NOT likely to contribute to the corresponding centroid. Given the new centroids, current learning scheme, where samples are enforced to learn from their assigned centroids solely, is unwise. To remedy the situation, we propose to use Confidence-Guided pseudo Label (CGL), which enables samples to approach not only the originally assigned centroid but other centroids that are potentially embedded with their identity information. Empowered by confidence-guided centroids and labels, our method yields comparable performance with, or even outperforms, state-of-the-art pseudo label refinement works that largely leverage auxiliary information.

1. Introduction

Person re-identification (ReID) aims to retrieve a person of interest across multiple cameras [14, 29, 35]. Due to the label-free training manner, unsupervised person ReID methods have attracted increasing attention. Unsupervised ReID methods can be broadly categorized into two types: unsupervised domain adaptation (UDA) methods [5, 9, 11, 25, 31, 40] and purely unsupervised learning (USL) methods [2, 4, 6, 23, 33]. The former pre-trains a model on person-related datasets, i.e., source domain, and fine-tunes it on ReID-related datasets, i.e., target domain. Apart from re-

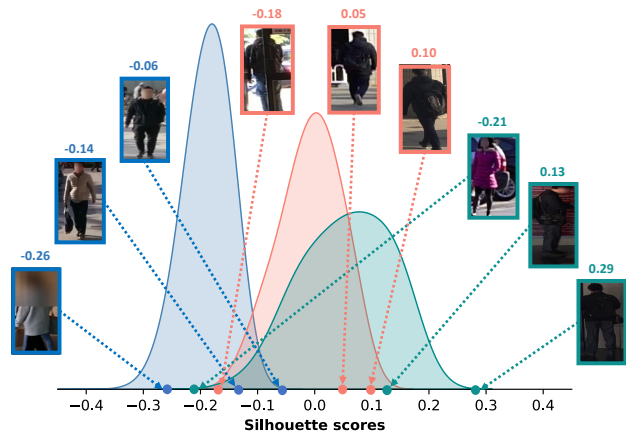


Figure 1. Training samples (cluster ID = 1) and their silhouette scores at epoch 0 (blue), epoch 25 (orange), and epoch 50 (green) on MSMT17 [24]. Higher silhouette scores denote samples are clustered at higher confidence. **Best viewed in color.**

quiring additional annotated labels, UDA methods are vulnerable to the large gap between the source domain and the target domain. In contrast, USL methods do not require any labeled data for training, which are more challenging but well fit real-world scenarios. In the paper, we focus on USL methods.

Existing USL methods generally follow a two-stage training scheme: 1) clustering, i.e., obtaining the pseudo labels via a clustering algorithm such as DBSCAN [8], and 2) network training, i.e., optimizing the network in a “supervised” manner with assigned cluster IDs. Contrastive loss such as InfoNCE [9] or ClusterNCE [6] usually serves as training objectives. Due to the blind trust in imperfect clustering results, the learning is inevitably misled by unreliable pseudo labels, where multiple identities are merged into one cluster or samples of one person are assigned to multiple clusters. Despite that some pseudo label refinement [2, 4, 23, 32, 33] have been proposed, they generally leverage auxiliary information, such as camera IDs [2, 23], body part predictions [4], and generated samples [33]. Given the fact that such auxiliary information is not free in reality, refining pseudo labels by merely exploiting internal characteristics within samples, i.e., the sample-

wise clustering confidence, appears to be more valuable.

To measure the sample-wise clustering confidence, *i.e.*, how well a sample fits its cluster, we employ a metric: silhouette score [19]. The score presents the ratio between intra-cluster distance and inter-cluster distance, which ranges from -1 to +1 (*higher is better*). To demonstrate the relationship between the clustering confidence and the silhouette score, we visualize silhouette scores of training samples of MSMT17 [24] in Fig. 1. Samples are from the same cluster (cluster ID=1) but at different training epochs, *i.e.*, 0, 25, and 50, respectively. As training goes on, the clustering is gradually enhanced by involving more effective features and the more discriminative network. At first images are grouped by coarse visual features, yet by identity-related information in the end. Meanwhile, sample-wise silhouette scores continuously shift towards higher values during training. Given this consistency, a conclusion can be drawn that, *a higher silhouette score implies the sample better fits its cluster, i.e., being clustered at higher confidence*. Previous learning schemes [6, 33] adopt all-sample based centroids, which are obtained by averaging features of all samples within the cluster, and enforce instances to approach such centroids. However, our observation suggests that low-confidence samples either are poor in quality or belong to other identities. Features of such images will inevitably contaminate centroids regardless of the training stage. In light of this, we propose Confidence-Guided Centroids (CGC) to provide more reliable cluster-wise prototypes for feature learning.

Although the reliability of cluster centroids has been improved, the conventional one-hot labeling strategy aggravates a problem. Since high-confidence samples exclusively contribute to the formation of cluster centroids, the identity-related information of low-confidence samples can hardly be presented in the assigned centroid. To illustrate the problem, an analysis is conducted on MSMT17 [24], where we intend to investigate how much identity information of low-confidence samples can be presented in their assigned centroids. We found that, with the vanilla all-sample based cluster centroids, only 5.83% low-confidence samples have their identity information embedded in the assigned centroid at the beginning. Although the ratio gradually climbs to 17.19%, a large proportion of low-confidence samples (over 80%) still are pushed to “wrong” centroids. Unfortunately, the ratio achieves 14.17% at most with confidence-guided centroids. Given the situation, the one-hot labeling strategy, which enforces samples to learn from the assigned centroid solely, is unwise. To address the problem, we propose to use confidence-guided pseudo labels (CGL), which encourages instances to approach not only the assigned confidence-guided centroid but also others where their identity information are potentially embedded.

In summary, our contributions are as follows:

- We propose Confidence-Guided Centroids (CGC) to provide cluster-wise prototypes for feature learning. The reliability of centroids is improved via filtering out low-confidence samples during formation.
- To overcome the problem that the identity information of low-confidence samples is rarely presented in their assigned centroids, we propose to use confidence-guided pseudo labels (CGL) during training. Apart from the originally assigned centroid, instances are also encouraged to approach other centroids where their identity information are potentially embedded.
- The proposed method only exploits internal characteristic for unsupervised person re-identification. Extensive experiments on benchmark datasets demonstrate that, our method yields better or comparable performances with state-of-the-art ones that largely leverage auxiliary information.

2. Related work

Unsupervised Person ReID. The existing unsupervised person ReID methods are divided into two categories: a) Unsupervised Domain Adaptation (UDA) methods, which boost the performance by leveraging the knowledge transferred from the source domain [5, 9, 11, 25, 31, 40], and b) purely UnSupervised Learning (USL) methods, which do not require any identity labels during training [2, 4, 6, 17, 32, 33]. Since UDA methods are highly prone to be affected by the large gap between the source domain and the target domain, they are hardly applicable to real-world scenarios [15, 29]. In the paper, we focus on USL methods.

Generally, USL methods exploit pseudo labels, instead of actual identity labels, as the guidance during training. Pseudo labels can be generated either by the image similarity [17, 22] or clustering algorithms [6, 16, 30, 33]. Specifically, SSL [17] and MMCT [22] formulate unsupervised person ReID as a classification task and predict pseudo labels based on the image similarity. In terms of clustering-based methods, BUC [16] and HCT [30] employ the bottom-up clustering scheme to gradually merge similar individual samples into clusters. Recently, Cluster-Contrast [6] adopts a contrastive learning scheme, which initializes, updates, and performs contrastive loss computation at the cluster level. However, clustering-based methods are generally sensitive to the pseudo label noise brought by imperfect clustering results.

Noise Reduction of Pseudo Label. Recently, how to handle noise pseudo labels in clustering-based methods has become a research hotspot. Specifically, SpCL [9] employs a self-paced learning scheme to gradually obtain more reliable clusters for the pseudo label refinement. CAP [23] splits each cluster into multiple proxies according to camera IDs. Such camera-aware proxies eliminate the pseudo label

noise brought by varying viewing points. ICE [2] alleviates the label noise by enhancing the consistency between augmented and original instances. RLCC [32] refines pseudo labels with the clustering consensus, which encourages the consistency between cluster results of two consecutive iterations. PPLR [4] employs the complementary relationship between reliable features of human global and body parts for the pseudo label refinement. ISE [33] generates boundary samples from actual samples and their neighboring clusters. The discriminability of the network is improved by enforcing generated samples to be correctly classified.

Unlike the above methods, this work explores whether internal characteristics can facilitate the pseudo label refinement. In the paper, we investigate the sample-wise clustering confidence, which describes how well a sample fits its cluster. With such criteria, we propose to employ better cluster centroids and pseudo labels for feature learning.

3. Methodology

3.1. Problem Statement

Let $\mathcal{T} = \{x_i\}_{i=1}^N$ denote an unlabeled training dataset, where x_i represents i -th image and N is the number of images. The USL ReID task aims to train a feature extractor E_θ in the unsupervised manner, where ReID features $\mathcal{F} = \{f_i\}_{i=1}^N$ are derived. The identity retrieval during inference is based on such ReID features. The training scheme of clustering-based USL methods [6, 9, 23, 33] alternates between two stages:

Stage I: Clustering. At the beginning of each epoch, training samples are clustered by DBSCAN [8]. Cluster IDs $y_i \in \{1, \dots, C\}$ serve as one-hot pseudo labels for the network optimization. Meanwhile, based on clustering results, a cluster-based memory bank $\mathcal{M} = \{m_i\}_{i=1}^C$ is initialized by cluster centroids that are formulated as,

$$m_i = \frac{1}{|\mathcal{C}|} \sum_{f_i \in \mathcal{C}} f_i, \quad (1)$$

where f_i represents the feature of i -th sample in the cluster \mathcal{C} , and $|\mathcal{C}|$ denotes the cluster size.

Stage II: Network Training. With the obtained pseudo labels, the network is then optimized in a ‘‘supervised’’ manner with the training objective, *i.e.*, ClusterNCE [6], which is formulated as,

$$\mathcal{L} = -\log \frac{\exp(\Phi(f \cdot m_+)/\tau)}{\sum_{j=1}^C \exp(\Phi(f \cdot m_j)/\tau)}, \quad (2)$$

where m_+ refers to the centroid of the cluster that f belongs to, m_j represents j -th centroid in the memory bank, $\Phi(u \cdot v)$ represents the cosine similarity between vector u and vector v , and τ is the temperature parameter. The memory bank is updated in a momentum manner [6] as,

$$m_i \leftarrow \mu \cdot m_i + (1 - \mu) \cdot f, \quad (3)$$

where μ is the updating factor and f refers to the feature of instance belonging to i -th cluster in the current mini-batch.

In this paper, we follow the framework of iterative clustering and network training. However, our method, as illustrated in Fig. 2, differs from previous works mainly in two aspects: 1) cluster centroids. Instead of using all samples to calculate the centroids, we adopt confidence-guided centroids (CGC) to provide reliable cluster-wise prototypes for feature learning (Sec. 3.3), and 2) pseudo labels. Apart from the assigned centroid, our confidence-guided pseudo labels (CGL) encourages instances to approach other centroids where their identity information are potentially embedded (Sec. 3.4).

3.2. Silhouette Score

To describe the sample-wise clustering confidence, *i.e.*, how well a sample fits its cluster, we employ a metric named silhouette score [19]. The score simultaneously considers two key factors of clustering, *i.e.*, tightness and separation.

Formally, for i -th data point in cluster \mathcal{C}_I , its average distance to other data points within the cluster can be calculated as,

$$a_i = \frac{1}{|\mathcal{C}_I|} \sum_{i,j \in \mathcal{C}_I, i \neq j} d(i, j), \quad (4)$$

where $d(i, j)$ refers to the distance between i -th and j -th data points and $|\mathcal{C}_I|$ represents the cluster size. Similarly, the distance between i -th data point and samples belonging to its nearest neighboring cluster \mathcal{C}_J can be denoted as,

$$b_i = \min_{J \neq I} \frac{1}{|\mathcal{C}_J|} \sum_{j \in \mathcal{C}_J} d(i, j). \quad (5)$$

Given the intra-class distance a_i and the minimal inter-class distance b_i , the silhouette score s_i is formulated as,

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}. \quad (6)$$

The silhouette score ranges from $[-1, 1]$. Note that, the score of clusters consisting of a single data point is 0. If an instance has a higher silhouette score, it has a smaller intra-class distance and a large inter-class distance. In other words, it is clustered at a higher confidence [19].

3.3. Confidence-guided Centroids

Based on the observation that images with lower silhouette scores (confidence) are generally containing high uncertainty regarding person identity, previous all-sample based cluster centroids are undoubtedly unwise. To remedy the problem, we build confidence-guided centroids (CGC) with high-confidence images only.

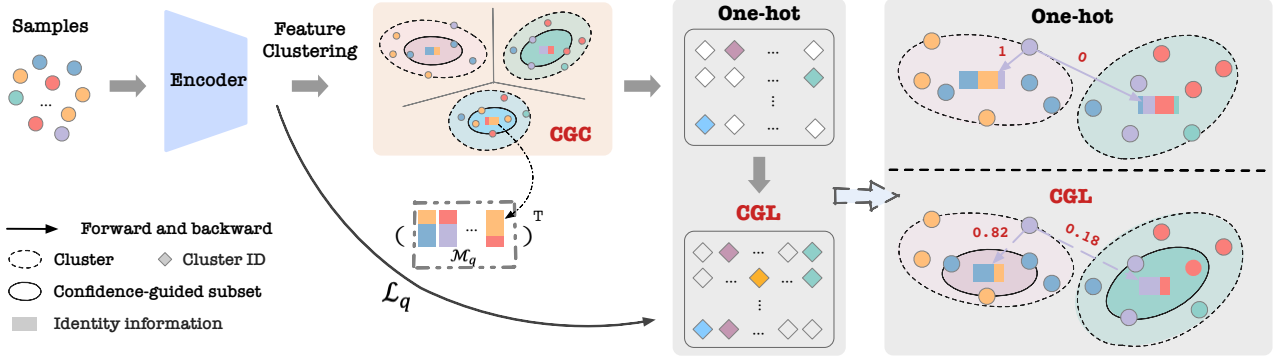


Figure 2. **Framework of the proposed method.** At the beginning of each epoch, training samples are clustered by DBSCAN [8]. Based on the original clustering result, we select a confidence-guided subset to build our confidence-guided centroids (CGC). During optimization, samples are encouraged to approach not only the assigned centroid but others where their identity information are potentially embedded via our confidence-guided pseudo labels (CGL).

Specifically, the confidence-guided centroid of i -th cluster m_i can be formulated as,

$$m_i = \frac{1}{|\mathcal{C}_q|} \sum_{f_i \in \mathcal{C}_q} f_i, \quad \mathcal{C}_q = \{f_i \in \mathcal{C} | s_i > \delta\}, \quad (7)$$

where a confidence-guided subset \mathcal{C}_q is selected from the original cluster \mathcal{C} by a silhouette score threshold δ . All confidence-guided centroids are then stored in a confidence-guided memory bank $\mathcal{M}_q = \{m_i\}_{i=1}^C$ for network optimization.

According to Fig. 1, our confidence-guided centroids can filter out images that are poor in quality or with cluttered backgrounds at early stages. While at later stages, such centroids effectively exclude some low-confidence samples that possibly belong to other identities. In summary, the proposed confidence-guided centroids can provide more reliable cluster-wise prototypes for feature learning.

3.4. Confidence-guided Pseudo Labels

Another problem of the clustering-based USL methods is that, samples, especially low-confidence ones, very likely carry different identity information with their assigned centroids. Our confidence-guided centroids also confronts with the problem since only high-confidence samples are included in the formation of centroids, as illustrated in Fig. 2. Given the situation, the previous learning scheme, which enforces samples to approach their assigned centroids solely regardless of the identity consistency in-between, is unwise. To alleviate the problem, we propose to use confidence-guided pseudo labels (CGL). Such labeling encourages samples to approach not only the assigned centroid but other centroids where their identity information are potentially embedded.

Specifically, we build a distance matrix $\mathcal{D} \in \mathbf{R}^{N \times C}$, where N and C denote the number of samples and clusters at the current epoch, respectively. In the paper, clus-

ters consisting of one sample are ignored [6]. As normalized identity features and centroids are adopted, $\mathcal{D}(i, j)$ represents the cosine distance between i -th sample and j -th confidence-guided centroid. Since similar samples are more likely to be scattered in neighboring clusters [33], the identity information of boundary samples is probably embedded in neighboring centroids. Therefore, when setting the learning target for samples, neighboring centroids should be assigned with higher confidence while distanced ones should be given lower confidence. To this end, a confidence matrix $\mathcal{P} \in \mathbf{R}^{N \times C}$ is obtained by,

$$\mathcal{P}(i, j) = \frac{p_{i,j}}{\sum_{j=1}^C p_{i,j}}, \quad p_{i,j} = \sigma(-\mathcal{D}(i, j)), \quad (8)$$

where $\mathcal{P}(i, j)$ represents the confidence of j -th centroid given by i -th sample, $\sum_{j=1}^C \mathcal{P}(i, j) = 1$, and $\sigma(\cdot)$ is the Sigmoid function. By integrating the confidence matrix with the originally assigned one-hot pseudo label y_i , the confidence-guided pseudo label of i -th sample \tilde{y}_i can be formulated as,

$$\tilde{y}_i = \beta \cdot y_i + (1 - \beta) \cdot \mathcal{P}(i), \quad (9)$$

where $\beta \in [0, 1]$ is the coefficient for the pseudo label refinement.

According to a previous work [26], the training objective, *i.e.*, ClusterNCE, can be considered as a non-parametric classifier, where centroids stored in the memory bank serve as the weight matrix of the classification layer. Therefore, the training objective of our method can be rewritten as,

$$\mathcal{L}_q = \frac{1}{N} \sum_{i=1}^N \left[\ell_{ce}(\mathcal{M}_q^T f_i, \tilde{y}_i) \right], \quad (10)$$

where ℓ_{ce} refers to the cross-entropy loss. Compared to Eq. (2), the training objective of our method can be obtained by simply applying two modifications: 1) replacing

Algorithm 1: Pipeline of our method

```
1 Require: Unlabeled data with pseudo labels
    $\mathcal{T} = \{(x_i, y_i)\}_{i=1}^N$ , where  $y_i \in \{1, \dots, C\}$ 
2 Require: Initialize the backbone encoder  $E_\theta$ 
3 Require: Threshold  $\delta$  for Eq. (7)
4 Require: Coefficient  $\beta$  for Eq. (9)
5 for  $n$  in  $[1, \text{epoch\_num}]$  do
6   Extracting features  $\mathcal{F}$  by  $E_\theta$ 
7   Clustering  $\mathcal{F}$  into  $C$  clusters with DBSCAN
8   Building CGC dictionary  $\mathcal{M}_q$  by Eq. (7)
9   for  $m$  in  $[1, \text{iteration\_num}]$  do
10    Sampling a mini-batch from  $\mathcal{T}$ 
11    Computing CGL with Eq. (9)
12    Computing loss with Eq. (10)
13    Updating encoder  $E_\theta$ 
14    Updating centroids with Eq. (3)
15  end
16 end
```

the original \mathcal{M} with our confidence-guided memory bank \mathcal{M}_q , and 2) replacing the one-hot pseudo label y_i with our confidence-guided one \tilde{y}_i . The training details are presented in Algorithm 1.

4. Experiment

4.1. Datasets and Evaluation Protocol

Datasets. We evaluate our proposed method on Market-1501 [34] and MSMT17 [24]. Market-1501 includes 32,668 images of 1,501 identities captured by 6 cameras. Among them, 12,936 images of 751 identities are used for training while the resting 19,732 images of 750 identities form the test set. MSMT17 contains 126,441 images from 4,101 identities captured by 15 cameras. The training set is composed of 32,621 images of 1,041 identities and the test set consists of 93,820 images of 3,060 identities. MSMT17 is more challenging due to the diversity in backgrounds, illuminations, poses, and occlusions.

Evaluation Protocol. Following previous methods [2, 6, 9, 33], the mean average precision (mAP) [1] and the cumulative matching characteristic (CMC) [34] top-1, top-5, top-10 accuracies are adopted as evaluation metrics. Note that, there are no post-processing operations, such as reranking [38], during inference.

4.2. Implementation Details

Following previous works [6, 9, 33], we adopt ResNet-50 [10] pre-trained on ImageNet [7] as our backbone feature encoder [6]. All layers after layer-4 are replaced by a generalized mean pooling (GeM) [18] layer followed by the batch normalization layer [12]. The output 2048-dimensional

ReID features are firstly normalized and then used for identity retrieval during inference. Our framework is built upon a state-of-the-art USL method [6]. For a fair comparison, we follow all experimental settings except for the formation of cluster centroids and the training objectives, as described in Sec. 3. The coefficient β in Eq. (9) is empirically set as 0.8 to achieve optimal performances.

During training, input images are resized to 256×128 . We adopt random flipping, cropping, and erasing [39] as data augmentation. Each mini-batch is formed by 16 identities, each with 16 images. Both identity and images are randomly selected from the training set. For the optimization, we adopt Adam [13] optimizer with a weight decay of 0.0005. The learning rate is set to 3.5×10^{-4} initially, and is divided by 10 every 30 epochs. We train for a total of 70 epochs on Market-1501 [34], and 50 on MSMT17 [24].

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

We compare our method with state-of-the-art (SOTA) unsupervised person ReID methods in Table 1. Compared with SOTA USL methods, our method outperforms previous ones, except ISE [33], on both benchmarks. Specifically, our method achieves 85.3% mAP and 94.2% top-1 accuracy on Market-1501 and 34.6% mAP and 63.4% top-1 accuracy on MSMT17. As stated in Sec. 2, existing SOTA methods generally leverage auxiliary information to refine pseudo labels. For example, CAP [23] and ICE [2] leverage the camera information, PPLR [4] employs body part predictions, and ISE [33] generates extra support samples in the latent space. As a departure from the above methods, our method yields SOTA performances by involving internal characteristics, *i.e.*, the sample-wise clustering confidence, only.

Additionally, we report the performance of some well-known supervised person ReID methods [20, 37] and unsupervised one [2] under the supervised setting in Table 1. Despite the absence of identity labels, our method even outperforms some supervised person ReID methods. Additionally, by replacing the pseudo labels with the ground-truth identity labels provided by datasets, our method outperforms an USL method (ICE [2]), which proves the potential of our framework.

4.4. Ablation Study

In this section, we thoroughly analyze the effectiveness of the proposed strategies, *i.e.*, confidence-guided centroids (CGC) and confidence-guided pseudo labels (CGL).

Effectiveness of CGC. We compare models trained with the vanilla all-sample based cluster centroids (“Baseline”) and with the proposed confidence-guided ones (“Baseline + CGC”). The performances are reported in Table 1. As can be seen, confidence-guided centroids boost the ReID performance by +1.7% / +0.6% on mAP / top-1 accuracy on

Method	Reference	Market-1501				MSMT17			
		mAP	top-1	top-5	top-10	mAP	top-1	top-5	top-10
<i>Purely Unsupervised</i>									
SSL [17]	CVPR'20	37.8	71.7	83.8	87.4	-	-	-	-
MMCL [22]	CVPR'20	45.5	80.3	89.4	92.3	11.2	35.4	44.8	49.8
HCT [30]	CVPR'20	56.4	80.0	91.6	95.2	-	-	-	-
SpCL [9]	NeurIPS'20	73.1	88.1	95.1	97.0	19.1	42.3	55.6	61.2
JNTL-MCSA [28]	CVPR'21	61.7	83.9	92.3	-	15.5	35.2	48.3	-
GCL [3]	CVPR'21	66.8	87.3	93.5	95.5	21.3	45.7	58.6	64.5
IICS [27]	CVPR'21	72.9	89.5	95.2	97.0	26.9	56.4	68.8	73.4
JVTC+* [3]	CVPR'21	75.4	90.5	96.2	97.1	29.7	54.4	68.2	74.2
OPLG-HCD [36]	ICCV'21	78.1	91.1	96.4	97.7	26.9	53.7	65.3	70.2
CAP [†] [23]	AAAI'21	79.2	91.4	96.3	97.7	36.9	67.4	78.0	81.4
ICE [2]	ICCV'21	79.5	92.0	97.0	98.1	29.8	59.0	71.7	77.0
ICE [†] [2]	ICCV'21	82.3	93.8	97.6	98.4	38.9	70.2	80.5	84.4
Cluster-Contrast [6]	Arxiv'21	82.1	92.3	96.7	97.9	27.6	56.0	66.8	71.5
PPLR [4]	CVPR'22	81.5	92.8	97.1	98.1	31.4	61.1	73.4	77.8
PPLR [†] [4]	CVPR'22	84.4	94.3	97.8	98.6	42.2	73.3	83.5	86.5
ISE [33]	CVPR'22	84.7	94.0	97.8	98.8	35.0	64.7	75.5	79.4
Cluster-Contrast (*Baseline)	Arxiv'21	82.4	92.5	96.9	98.0	31.4	61.2	72.5	76.9
Baseline+CGC	-	84.1	93.1	97.2	98.2	34.1	63.1	75.0	79.0
Baseline+CGL	-	83.4	93.2	97.1	98.2	33.7	62.5	73.9	78.4
Ours	-	85.3	94.2	97.6	98.5	34.6	63.4	74.6	79.3
<i>Supervised</i>									
PCB [20]	ECCV'18	81.6	93.8	97.5	98.5	40.4	68.2	-	-
DG-Net [37]	CVPR'19	86.0	94.8	-	-	52.3	77.2	-	-
ICE (w/ ground-truth) [2]	ICCV'21	86.6	95.1	98.3	98.9	50.4	76.4	86.6	90.0
Our (w/ ground-truth)	-	87.4	95.3	98.5	99.0	51.0	76.6	87.1	90.1

Table 1. Comparison of ReID methods on Market-1501 and MSMT17 datasets. The best USL results without camera information are marked with **bold**. † indicates using the additional camera knowledge.

Market-1501, and +2.7% / +1.9% on MSMT17. Such improvements reveal the potential of the clustering confidence in the pseudo label refinement.

To better understand how our confidence-guided centroids benefit feature learning, we analyze how the sample-wise confidence varies throughout the training process on MSMT17. Specifically, we visualize the distribution of silhouette scores at different epochs in Fig. 3. Note that scores of outliers are excluded. Several conclusions can be drawn from the comparison between Fig. 3(a) and Fig. 3(b). 1) As training goes on, the number of valid samples gradually increases, representing as larger areas under the curve. 2) Starting from the same point (epoch 0), with our confidence-guided centroids, a noticeable shift towards higher scores can be found at epoch 25. The shift implies CGC can effectively reduce the overall number of low-confidence samples while enhancing high-confidence ones. 3) The advantage remains until the end of training. At epoch 50, the number of high-confidence samples increases, representing by a higher peak closer to 0.4.

Effectiveness of CGL. We also compare the baseline model (“Baseline”) and the model trained with confidence-guided pseudo labels (“Baseline + CGL”). The performances are shown in Table 1. As can be seen, CGL improves mAP and top-1 accuracy by 1.0% and 0.7% on Market-1501, by 2.3% and 1.3% on MSMT17. When both CGC and CGL are em-

ployed during training, improvements are +2.9% and +1.7% on Market-1501, and +3.2% and +2.2% on MSMT17.

In terms of the sample-wise clustering confidence, we visualize the distribution of silhouette scores in Fig. 3(c), when CGL is applied during training. Compared to the model trained without CGL (Fig. 3(b)), CGL further pushes the score towards a higher value at both epoch 25 and epoch 50. Less low-confidence samples during training implies our CGL contributes to better clustering. In summary, the above qualitative and quantitative results prove the proposed scheme can boost performance by enhancing the sample-wise clustering confidence.

4.5. Parameter Analysis

Threshold δ in CGC. To obtain the optimal threshold δ in Eq. (7) for the proposed confidence-guided centroids (CGC), three types of threshold selection strategies are explored, *i.e.*, linear, dynamic and constant, respectively. For the former two strategies, the threshold gradually increases as training goes on. The constant strategy employs a fixed threshold throughout the training process.

Specifically, the linear strategy updates the threshold by $\delta_t = \delta_0 * t / T + \epsilon$, where δ_0 limits the range of threshold and ϵ is the offset. In the paper, we set $\delta_0 = 0.2$ and $\epsilon = -0.1$. t and T denote the current epoch and the overall number of epochs, respectively. In terms of the dynamic strategy, the

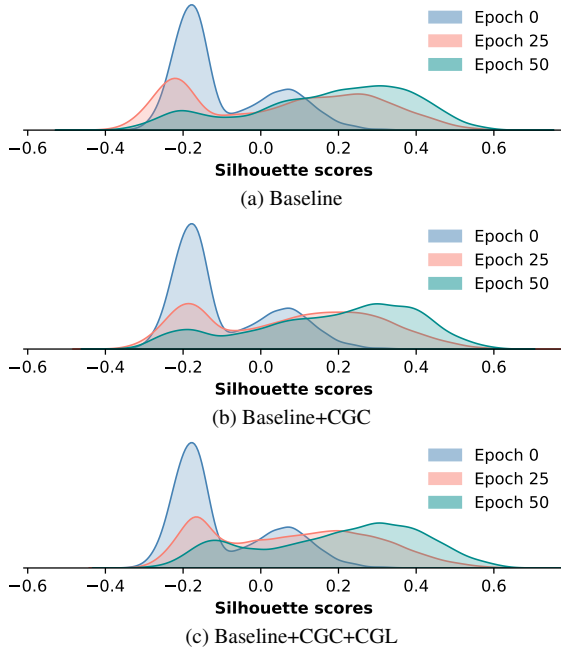


Figure 3. Silhouette scores of valid samples (MSMT17 [24]) at different epochs. Comparisons are conducted between (a) baseline model, (b) baseline model with confidence-guided centroids (CGC), and (c) baseline model with CGC and confidence-guided pseudo labels (CGL). **Best viewed in color.**

Method	Strategy	δ	Market-1501		MSMT17	
			mAP	top-1	mAP	top-1
Baseline	-	-	82.4	92.5	31.4	61.2
Ours	Linear	-	85.3	94.2	33.6	63.0
		-	84.9	93.9	33.0	62.8
	Constant	-0.1	83.5	93.4	32.7	62.8
		0	84.9	94.0	34.6	63.4
		0.1	84.0	93.3	34.0	63.2

Table 2. Comparison of threshold selection strategies of confidence-guided centroids (CGC) on benchmark datasets.

threshold is updated by $\delta = \delta_0 * \tanh(0.1 * (t - T/2))$. We set $\delta_0 = 0.1$ to achieve $\delta \in [-0.1, 0.1]$, which is the same as the linear strategy. The range is set empirically in the consideration of the image quality and the distribution of silhouette scores (see Fig. 3). Apart from the varying threshold, we conduct the constant strategy by fixing the threshold as $\{-0.1, 0, 0.1\}$ respectively. The comparisons between model performances with different strategies are reported in Table 2. The best performance is achieved when adopting the linear strategy for Market-1501 and applying a fixed threshold $\delta = 0$ on MSMT17. The optimal settings are employed in all experiments.

Coefficient β in CGL. To analyze the impact of the coefficient β in the proposed confidence-guided pseudo labels (CGL), we tune the value of parameter β from 0 to 1 while keeping others fixed. According to Eq. (9), when β is set

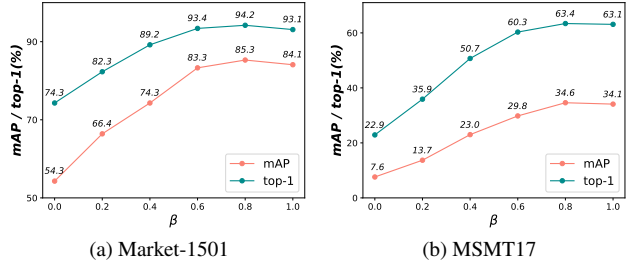


Figure 4. Comparison of coefficient β in confidence-guided pseudo labels (CGL) on (a) Market-1501 and (b) MSMT17.

to 0 or 1, our method decomposes down to using the confidence matrix or the one-hot pseudo label exclusively during training. The results on two benchmarks are illustrated in Fig. 4. As shown, as β increases from 0 to 0.8, both mAP and top-1 accuracy increase. A slight performance drop can be found when increasing β from 0.8 to 1. To achieve the best performance, we set $\beta = 0.8$ for all experiments.

4.6. More Discussions

Identity Feature Distribution. To better understand the advantages of the proposed strategies, we visualize the distribution of identity features via t-SNE [21]. Specifically, 20 identities are randomly selected from Market-1501 [34] and MSMT17 [24], respectively. Features of selected identities are extracted by the baseline model and our model is trained with confidence-guided centroids (CGC) and confidence-guided pseudo labels (CGL). The distribution of identity features is illustrated in Fig. 5. As can be seen, due to the vast variety in camera views, backgrounds, and poses, the feature distribution of MSMT17 is more chaotic than that of Market-1501. Despite such challenges, with the aid of the proposed strategies, features of the same identity are distributed more compactly while features of different identities are further separated.

Identity Consistency Score. The current learning scheme enforces samples to approach their assigned cluster centroids, where their identity information are embedded. However, the existence of noisy labels will lead samples to “wrong” centroids. It is especially problematic for low-confidence samples, *i.e.*, boundary samples, because they can be closer to other centroids than the assigned ones.

To investigate the problem, we conduct an experiment on MSMT17 to analyze how much the identity information of boundary samples can be presented in the assigned centroids, *i.e.*, the identity consistency in-between. Specifically, we select clusters whose size is over 100 at each epoch. For each cluster, samples whose silhouette scores rank at the bottom 5% are empirically marked as boundary samples. Formally, let $\mathcal{C} = \{(x_i, g_i)\}_{i=1}^{N_c}$ denote a cluster with N_c samples, where g_i refers to the ground-truth identity label provided by the dataset. An identity

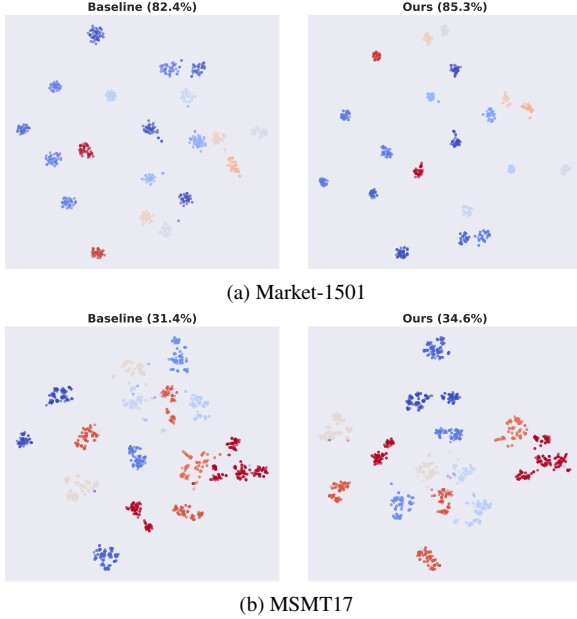


Figure 5. Visualization of the identity feature distribution via t-SNE [21] on (a) Market-1501 and (b) MSMT17. For each group, features are derived by the baseline model (left) and the model trained with the proposed confidence-guided centroids (CGC) and pseudo labels (CGL) (right), respectively. Model performances (mAP) are also denoted. Different identities are denoted by different colors. **Best viewed in color.**

set $\mathcal{G} = \{g_k\}_{k=1}^M$ is then constructed by overall M identities occurring in the cluster. Following the formation of vanilla all-sample based cluster centroids (Eq. (1)), the identity information embedded in the centroid can be obtained by linearly integrating all identities within the cluster via weights $\mathcal{Q} = \{q_k\}_{k=1}^M$, where q_k is obtained by $q_k = \frac{1}{|\mathcal{C}|} \sum_{g_i \in \mathcal{C}} \mathbb{1}\{g_i = g_k\}$. $|\mathcal{C}|$ denotes the cluster size. $\mathbb{1}\{g_i = g_k\}$ equals to 1 when $g_i = g_k$, otherwise 0. Then, the identity consistency score (ICS) between boundary samples and the cluster centroid of \mathcal{C} can be calculated as, $ICS = \frac{1}{N_c} \sum_{g_i \in \mathcal{C}} q_k \cdot \mathbb{1}\{g_i = g_k\}$.

Similar to the vanilla scheme, ICS of our confidence-guided centroids (CGC) scheme can be computed by simply replacing \mathcal{C} with the confidence-guided subset \mathcal{C}_q during the computation of the weight q_k . Since low-confidence samples are filtered out in the formation of confidence-guided centroids, the identity set \mathcal{G} only includes identities of samples with high confidence scores. We compare the average ICS throughout the training with vanilla all-sample based cluster centroids and the proposed confidence-guided ones, and obtain the curves in Fig. 6.

For the vanilla scheme, only 5.83% boundary samples carry the same identity information with their assigned centroid at the beginning. Although the ratio gradually climbs to 17.19%, a large proportion of boundary samples (over 80%) still are pushed to centroids where their identity in-

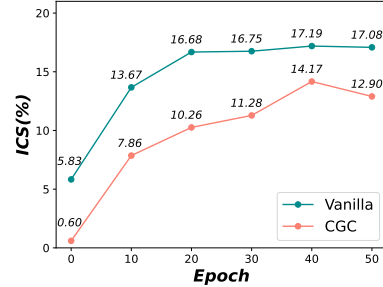


Figure 6. Identity consistent score (ICS) of boundary samples at different epochs. Vanilla and CGC refer to the previous all-sample based cluster centroids and the proposed confidence-guided centroids, respectively.

formation are rarely presented. Unfortunately, the problem has been aggravated by confidence-guided centroids, where the ratio achieves 14.17% at most. The low identity consistency scores point out the seriousness of the problem and validates the necessity of our confidence-guided pseudo labels.

5. Conclusion

This paper focused on the pseudo label refinement for clustering-based unsupervised person ReID, which aims to alleviate the pseudo label noise brought by imperfect clustering results. Instead of relying on auxiliary information such as camera IDs, body parts, or generated samples, we refined pseudo labels with internal characteristics, *i.e.*, the sample-wise clustering confidence. Specifically, we proposed to use confidence-guided centroids (CGC) to provide reliable cluster-wise prototypes for feature learning, where low-confidence instances are filtered out during the formation of centroids. Additionally, targeting at the problem that a large proportion of samples are pushed to “wrong” centroids, we propose to use confidence-guided pseudo labels (CGL). Such labeling enables samples to approach not only the assigned centroid but other clusters where their identities are potentially embedded. With the aid of CGC and CGL, our method yields comparable performances with, or even outperforms, state-of-the-art pseudo label refinement works that largely leverage auxiliary information.

Limitations and Broader Impact. Although we conducted multiple threshold strategies in the paper, the range is selected empirically. We are interested in exploring adaptive thresholds in the future. Despite that our method did NOT leverage either identity labels or auxiliary information, it may still involve a concern for human privacy during the data collection. Therefore, the legal utilization of person ReID data should be regulated strictly to avoid ethical issues.

References

- [1] Song Bai, Xiang Bai, and Qi Tian. Scalable person re-identification on supervised smoothed manifold. In *CVPR*, 2017. 5
- [2] Hao Chen, Benoit Lagadec, and Francois Bremond. Ice: Inter-instance contrastive encoding for unsupervised person re-identification. In *ICCV*, 2021. 1, 2, 3, 5, 6
- [3] Hao Chen, Yaohui Wang, Benoit Lagadec, Antitza Dantcheva, and Francois Bremond. Joint generative and contrastive learning for unsupervised person re-identification. In *CVPR*, 2021. 6
- [4] Yoonki Cho, Woo Jae Kim, Seunghoon Hong, and Sung-Eui Yoon. Part-based pseudo label refinement for unsupervised person re-identification. In *CVPR*, 2022. 1, 2, 3, 5, 6
- [5] Yongxing Dai, Jun Liu, Yifan Sun, Zekun Tong, Chi Zhang, and Ling-Yu Duan. Idm: An intermediate domain module for domain adaptive person re-id. In *ICCV*, 2021. 1, 2
- [6] Zuozhuo Dai, Guangyuan Wang, Weihao Yuan, Siyu Zhu, and Ping Tan. Cluster contrast for unsupervised person re-identification. *arXiv:2103.11568*, 2021. 1, 2, 3, 4, 5, 6
- [7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 5
- [8] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *KDD*, 1996. 1, 3, 4
- [9] Yixiao Ge, Feng Zhu, Dapeng Chen, Rui Zhao, and Hongsheng Li. Self-paced contrastive learning with hybrid memory for domain adaptive object re-id. In *NeurIPS*, 2020. 1, 2, 3, 5, 6
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5
- [11] Tao He, Leqi Shen, Yuchen Guo, Guiguang Ding, and Zhenhua Guo. Secret: Self-consistent pseudo label refinement for unsupervised domain adaptive person re-identification. In *AAAI*, 2022. 1, 2
- [12] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 5
- [13] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5
- [14] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *CVPR*, 2018. 1
- [15] Xiangtan Lin, Pengzhen Ren, Chung-Hsing Yeh, Lina Yao, Andy Song, and Xiaojun Chang. Unsupervised person re-identification: A systematic survey of challenges and solutions. *arXiv:2109.06057*, 2021. 2
- [16] Yutian Lin, Xuanyi Dong, Liang Zheng, Yan Yan, and Yi Yang. A bottom-up clustering approach to unsupervised person re-identification. In *AAAI*, 2019. 2
- [17] Yutian Lin, Lingxi Xie, Yu Wu, Chenggang Yan, and Qi Tian. Unsupervised person re-identification via softened similarity learning. In *CVPR*, 2020. 2, 6
- [18] Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human annotation. *IEEE TPAMI*, 2018. 5
- [19] Peter J Rousseeuw. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. *Comput. Appl. Math.*, 1987. 2, 3
- [20] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *ECCV*, 2018. 5, 6
- [21] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 2008. 7, 8
- [22] Dongkai Wang and Shiliang Zhang. Unsupervised person re-identification via multi-label classification. In *CVPR*, 2020. 2, 6
- [23] Menglin Wang, Baisheng Lai, Jianqiang Huang, Xiaojin Gong, and Xian-Sheng Hua. Camera-aware proxies for unsupervised person re-identification. In *AAAI*, 2021. 1, 2, 3, 5, 6
- [24] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *CVPR*, 2018. 1, 2, 5, 7
- [25] Yuhang Wu, Tengting Huang, Haotian Yao, Chi Zhang, Yuanjie Shao, Chuchu Han, Changxin Gao, and Nong Sang. Multi-centroid representation network for domain adaptive person re-id. In *AAAI*, 2022. 1, 2
- [26] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *CVPR*, 2018. 4
- [27] Shiyu Xuan and Shiliang Zhang. Intra-inter camera similarity for unsupervised person re-identification. In *CVPR*, 2021. 6
- [28] Fengxiang Yang, Zhun Zhong, Zhiming Luo, Yuanzheng Cai, Yaojin Lin, Shaozi Li, and Nicu Sebe. Joint noise-tolerant learning and meta camera shift adaptation for unsupervised person re-identification. In *CVPR*, 2021. 6
- [29] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. Deep learning for person re-identification: A survey and outlook. *IEEE TPAMI*, 2021. 1, 2
- [30] Kaiwei Zeng, Munan Ning, Yaohua Wang, and Yang Guo. Hierarchical clustering with hard-batch triplet loss for person re-identification. In *CVPR*, 2020. 2, 6
- [31] Xinyu Zhang, Jiewei Cao, Chunhua Shen, and Mingyu You. Self-training with progressive augmentation for unsupervised cross-domain person re-identification. In *ICCV*, 2019. 1, 2
- [32] Xiao Zhang, Yixiao Ge, Yu Qiao, and Hongsheng Li. Refining pseudo labels with clustering consensus over generations for unsupervised object re-identification. In *CVPR*, 2021. 1, 2, 3
- [33] Xinyu Zhang, Dongdong Li, Zhigang Wang, Jian Wang, Er-rui Ding, Javen Qinfeng Shi, Zhaoxiang Zhang, and Jingdong Wang. Implicit sample extension for unsupervised person re-identification. In *CVPR*, 2022. 1, 2, 3, 4, 5, 6
- [34] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *ICCV*, 2015. 5, 7

- [35] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. Person re-identification in the wild. In *CVPR*, 2017. [1](#)
- [36] Yi Zheng, Shixiang Tang, Guolong Teng, Yixiao Ge, Kaijian Liu, Jing Qin, Donglian Qi, and Dapeng Chen. Online pseudo label generation by hierarchical cluster dynamics for adaptive person re-identification. In *ICCV*, 2021. [6](#)
- [37] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *CVPR*, 2019. [5](#), [6](#)
- [38] Zhun Zhong, Liang Zheng, Donglin Cao, and Shaozi Li. Re-ranking person re-identification with k-reciprocal encoding. In *CVPR*, 2017. [5](#)
- [39] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data augmentation. In *AAAI*, 2020. [5](#)
- [40] Zhun Zhong, Liang Zheng, Zhiming Luo, Shaozi Li, and Yi Yang. Invariance matters: Exemplar memory for domain adaptive person re-identification. In *CVPR*, 2019. [1](#), [2](#)