# Contrastive Multiple Instance Learning for Weakly Supervised Person ReID

Jacob Tyo
jtyo@cs.cmu.edu
DEVCOM Army Research Lab
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

Zachary C. Lipton
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

## ABSTRACT

The acquisition of large-scale, precisely labeled datasets for person re-identification (ReID) poses a significant challenge. Weakly supervised ReID has begun to address this issue, although its performance lags behind fully supervised methods. In response, we introduce Contrastive Multiple Instance Learning (CMIL), a novel framework tailored for more effective weakly supervised ReID. CMIL distinguishes itself by requiring only a single model and no pseudo labels, while leveraging contrastive losses – a technique that has significantly enhanced traditional ReID performance yet is absent in all prior MIL-based approaches. Through extensive experiments and analysis across three datasets, CMIL not only matches state-of-the-art performance on the large-scale SYSU-30k dataset with fewer assumptions but also consistently outperforms all baselines on the WL-market1501 and Weakly Labeled MUddy racer re-iDentification dataset (WL-MUDD) datasets. We introduce and release the WL-MUDD dataset, an extension of the MUDD dataset featuring naturally occurring weak labels from the real-world application at PerformancePhoto.co. All our code and data are accessible at https://drive.google.com/file/d/1rjMbWB6m-apHF3Wg_cfqc8QqKgQ21AsT/view?usp=drive_link.

## CCS CONCEPTS

• **Computing methodologies → Semi-supervised learning settings**; **Visual content-based indexing and retrieval**.

## KEYWORDS

Multiple Instance Learning, Weak Supervision, Contrastive Learning, Person Re-Identification, Machine Learning

## 1 INTRODUCTION

Accurate data labeling is a critical part of any machine-learning system, but is often prohibitively expensive, especially for person

re-identification (ReID). In most classification problems, the classes are easily human-recognizable, allowing annotators to quickly recognize and label the class of a data point. In ReID however, the data points can consist of millions of individuals, none of which are known to the annotators. In this case, generating accurate labels is extremely difficult and time-consuming. An alternative approach is to use weakly supervised learning (WSL) methods that can effectively leverage lower-quality data labeling, which is often available in larger amounts at meager cost [26, 29, 39, 42, 54].

WSL has achieved impressive results on benchmark datasets, but performance still lags that of the standard, fully-supervised, setting. Given that the reid task of identifying images of the same person is inherently contrastive with respect to identities, it seems possible that we could leverage techniques from contrastive learning to improve WSL further. Contrastive learning is a specific subset of supervised learning where models are optimized on pairs (or triplets) of inputs to determine if the inputs originate from the same class or not. This slight reframing has major benefits both in terms of model performance and generalization [12, 18], and in terms of computational efficiency in downstream applications. Specifically for ReID, common applications include facial recognition, person search, and image retrieval. In each of these settings, the number of downstream classes (identities) is typically unknown, a setting where contrastive models excel. However, traditional algorithms in contrastive learning depend on accurate labels.

Weak labels for ReID can be gathered in several ways. One example, as provided by Guillaumin et al. [14], is to gather images of people based on an online search. The resulting dataset is bags of images that all contain the same person, but the images would also be extremely noisy, containing many other people in each photo. Another example of this type of weak labels is to observe event photo purchases - someone purchasing photos of a racer after a marathon likely purchase photos that all contain a common single person. As part of this work, we introduce and release the Weakly Labeled MUddy racer re-iDentification dataset (WL-MUDD) dataset, which is a dataset labeled in this exact manner from the motorcycle racing event photo website PerformancePhoto.co.

The dominant methods for weak ReID rely on pseudo-labeling [26, 39, 54, 58], which is an iterative process of predicting new labels for the weakly labeled data in an attempt to build better models. Other approaches include graph-based methods [30, 42], Multiple Instance Learning (MIL) [19, 32, 43], or transferring an unsupervised model (i.e. trained without labels). The unsupervised methods have made significant progress recently, but still fall short of methods that can leverage labeling [39]. Within WSL, pseudo-labeling approaches typically outperform those of noisy learning and MIL. However, the

existing MIL formulations restrict the use of contrastive methodologies.

In this work, we introduce Contrastive Multiple Instance Learning to enable contrastive learning among weakly labeled bags of images. Contrastive learning is typically interpreted as decreasing the distance between the representation of two images of the same identity (or class), and increasing the distance between the representation of two images of different identities. However, in weakly supervised learning, the labels are not that granular. Pseudo-labeling methods get around this by trusting that the labels are granular enough, and then updates the labels as training progresses, but this is prone to errors. Especially in settings where the intra-identity variability is extremely high, and the inter-identity variability is low, which is the exact case for our WL-MUDD dataset. Instead of a label refinement approach, we focus on the MIL formulation, and to enable contrastive techniques, we formulate the contrastive learning problem as decreasing the distance between two *bag* representations that have the same label, and increase the distance between two *bag* representations with a different label.

This shift in perspective, of optimizing for bag representations instead of representations of a single image within that bag, is not obvious, mainly because at test time, the goal is still to produce a high-performing ReID model: one that can take a single image and produce a high-quality embedding for it. To this end, CMIL includes two processes to help in this regard. The first is that each image in each bag is independently embedded into a representation using a feature extraction network, resulting in a bag of image features. Then, the bag of image features is passed through an accumulation network to generate a bag representation. Second, we experiment with an *alignment loss*, to encourage our model to learn image and bag representations that are similar.

The feature extraction network is chosen to be a standard ReID model, specifically ResNet-50 [16]. The accumulation network must be permutation invariant, and is therefore chosen to be a set transformer, although we do provide ablation studies with the simpler choices of the average, max, and sum operators. Surprisingly, we find that even without the alignment loss, optimizing for high-quality bag representations implicitly leads to high-quality image representations.

We evaluate CMIL against a state-of-the-art weakly supervised learning method [49] and a prior MIL method [29] on the weakly-labeled Market1501 (WL-Market1501) dataset, and WL-MUDD datasets. The WL-Market1501 dataset is the widely used Market1501 dataset [55] but with noise added to mimic the weakly labeled setting. Then, we compare CMIL to the state-of-the-art on the large-scale SYSU-30k weakly labeled ReID dataset, containing nearly 30 million images and over 30 thousand identities. We find that on both WL-Market1501 and WL-MUDD, CMIL consistently achieved the best rank-1, rank-5, rank-10 accuracy, and mean accuracy precision. On the SYSU-30k dataset, CMIL matched the state-of-the-art while requiring fewer modeling assumptions. Lastly, our ablation studies reveal the surprising effectiveness of average pooling for image aggregation, along with the surprisingly different instance and bag representations even of the best-performing models.

The contributions of this work are threefold:

- The introduction and release of WL-MUDD, a real-world dataset of motorcycle racers with naturally weak labels from PerformancePhoto.co.
- We introduce CMIL, a novel framework for re-identification from weakly labeled group images.
- Experimental evidence of the efficacy of CMIL and an analysis highlighting the surprising differences between image and bag representations.

## 2 DATASETS AND PROBLEM SETUP

In this section, we formally introduce the weakly supervised ReID setting, as well as a new weakly supervised ReID dataset. The dataset is available at https://drive.google.com/file/d/1rjMbWB6m-apHF3Wg_cfqc8QqKgQ21AsT/view?usp=drive_link.

### 2.1 Weakly Supervised Re-Identification

The problem we address in this paper is re-identification from weakly labeled group images. We assume that we are given a dataset of images $I = \{I_1, I_2, \ldots, I_N\}$ where each image $I_j$ contains one or more people that we are interested in identifying. Let $X_j = \{x_1^j, x_2^j, \ldots, x_{M_j}^j\}$ denote the set of $M_j$ *crops* containing each person extracted from image $I_j$. We refer to each specific person in an image $I_j$ as $x_i^j \in X_j$.

However, unlike conventional re-identification datasets, we only have weak labels for each image. These labels merely indicate the presence of a shared identity within each group, but not the specific identity of each individual instance within the group. This means that we have access to a set of *bags* $\mathcal{B} = \{B_1, B_2, \ldots, B_K\}$ where each bag $B_k = \{I_{k_1}, I_{k_2}, \ldots, I_{k_{|B_k|}}\}$ contains images that share a common identity. Importantly, the individual instances within each group are not labeled with their specific identities. Instead, the bag is labeled with only a single identity. The key challenge in this setting is to learn a model that can effectively discriminate between different identities despite only having access to these weak bag-level labels. During inference time the bag-level labels are not the label of interest. Instead, we want a standard ReID model at inference time, meaning that we need to be able to predict the identity of a single person (i.e. crop) within an image, not of the bag.

### 2.2 WL-MUDD Dataset

PerformancePhoto.co is an online marketplace for off-road racing photographers and fans. Powered by text spotting and ReID models to enable searchable racing photos, improvements to the ReID models suffer from the high costs of ReID dataset labeling. However, there is a proxy that gives natural weak labels: user purchases. When a user purchases photos from a single event, they are likely purchasing photos of a single individual. However, it is also likely that there is more than one individual in each photo purchased. Following notation from Section 2.1, the set of photos purchased by a single user can be regarded as a *bag* that can be weakly labeled with a unique identity.

The MUDD [33] dataset was curated from PerformancePhoto.co and manually labeled in the traditional, fully supervised, ReID setting. We adapt the MUDD dataset to the weakly supervised
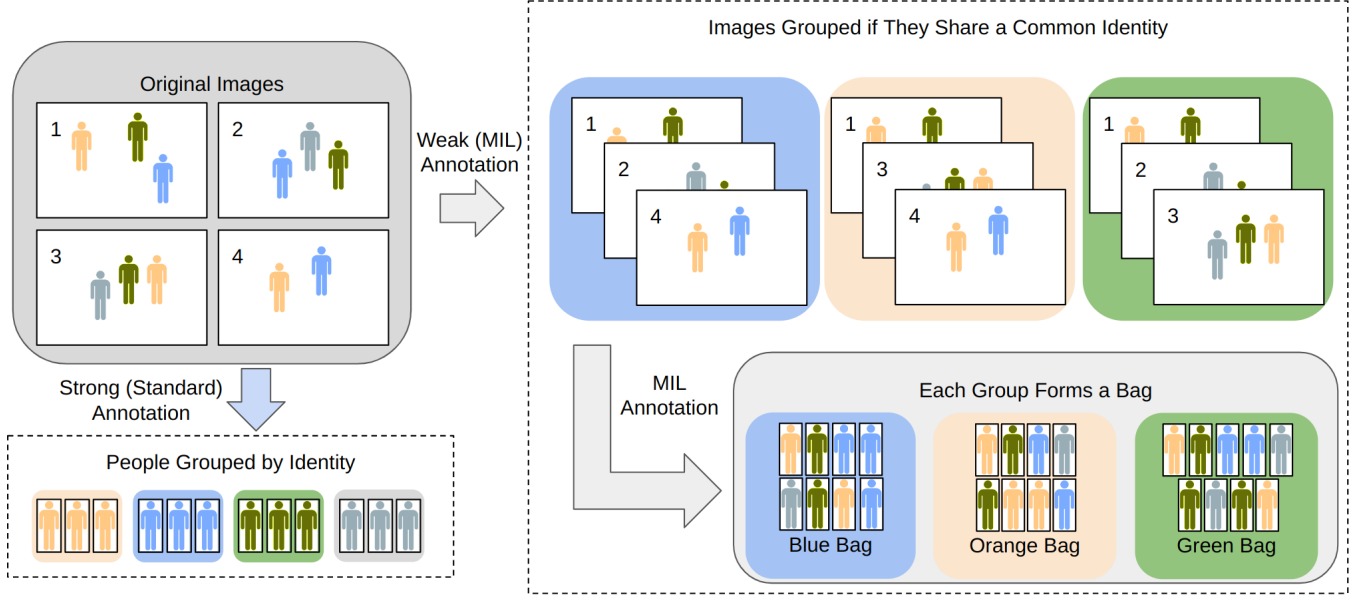
**Figure 1: The annotation process for strong and weak ReID. The strong annotations group each crop into a bag based on their identity, whereas the weak annotation groups all images based on a shared identity, and then all crops from the grouped images become a bag.**

setting by re-labeling the data points at the bag level and adding all, previously unlabeled, crops to each bag according to their existence in the original images. Figure 1 shows this labeling process and compares it to the standard (strong) annotation procedure. Instead of relying on the user purchases heuristic, we were able to build out the WL-MUDD dataset by taking all of the strong labels from the MUDD dataset, and then linking them back to the photos they originated from. Then, we take all the other people in the original photo, and add them to the dataset under the same label, forming a bag. This is repeated for every person in the original MUDD dataset, resulting in a weakly labeled dataset over twice as large. We refer to this dataset as the Weakly Labeled MUddy racer re-iDentification dataset (WL-MUDD).

The average bag in WL-MUDD has 75 crops of people in it, with 32% of them being the identity of the label attributed to that bag. This corresponds to an average noise level of 68%. The bags can be as small as 5 crops, or as large as 300, and the noise level of each bag varies between 50% and 85%. Figure 2 gives examples of bags in the dataset, highlighting the extremely high inter-class variation. The crops highlighted in green are representative of the bag label, whereas the red highlighted crops are not.

## 3 CONTRASTIVE MULTIPLE INSTANCE LEARNING

We cast the weakly supervised object re-identification problem as one of multiple-instance learning and present the contrastive multiple-instance learning (CMIL) method. A standard multiple-instance learning problem handles bag-level labeling by getting a feature representation for all crops in a bag, applying an accumulation function (typically max, average, etc.) to get a single bag

representation from all of the crop representations, and then applying a classifier to the bag representation to determine a classification. Instead of a bag classifier (or alongside), we compare bag representations via a contrastive loss. This allows us to train end-to-end in a contrastive fashion. The CMIL framework is shown in Figure 3.

This is a divergence from standard contrastive learning. At test time we compare representations of crops, and therefore the goal of training is to optimize the crop representations accordingly. However, in this formulation, we are directly optimizing the bag representations, and only indirectly optimizing the crop representations. Specifically, all crops from a single bag $k$ are encoded into crop representations by a model $f$ parameterized by $\theta$.

$$z_j^i = f_\theta(x_j^i), \ \forall j \in M_i, i \in N_k, \tag{1}$$

where $M_i$ is the number of crops in image $i$ and $N_k$ is the set of images in bag $k$. It is critical that this model takes a specific crop as input, and returns the corresponding representation for that input because during testing, this is the only aspect of the model that will be utilized. Then given all crop representations for a bag $k$, they must be accumulated into a single bag-representation using a model $g$ parameterized by $\phi$.

$$r_k = g_\phi(z_1^1, \ldots, z_{M_i}^{N_k}). \tag{2}$$

This accumulation function should be permutation invariant to the input, as there is no way to control the ordering of the instances meaningfully.

The final component of this architecture is a distance or similarity function $d$. To apply contrastive learning, we must be able to measure the distances between pairs/triplets/quadruplets of bags.

Figure 2: Four example subsets from four different bags of the WL-MUDD dataset. Each image within a bag is outlined in green if it is the same identity as the bag, and red if it is not. Each bag can have very different ratios of correct to incorrect identities of the underlying images.
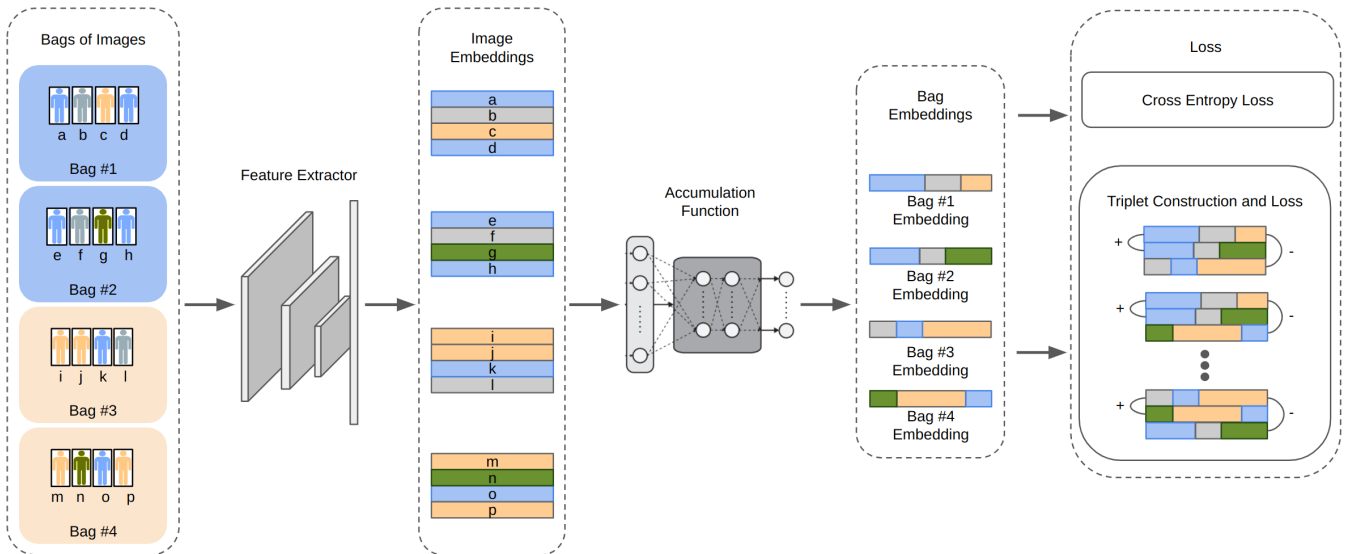


Figure 3: The CMIL framework. For each image in a batch of bags, a feature extraction network is used to get an embedding for each image. Then for each bag, the corresponding image embeddings are combined into a single bag embedding via an accumulation function. Finally, the bag embeddings are used to calculate the cross entropy loss (or identity loss), as well as the triplet loss based on all valid triplets from the batch.

Any proper distance metric, such as the Euclidean or cosine distance, can be used. This distance can then be used to return a ranking, thresholded to provide a classification, etc. We focus on the setting where we are given a triplet of bags (or by the time it reaches the distance metric, bag representations). Given a bag $a$ and $b$, the distance between their representations is represented by:

$$\hat{y} = d(r_a, r_b). \tag{3}$$

Note that in this methodology, we depend on bags of data during each iteration. Each iteration must have a sufficient number of bags, as well as a sufficient number of crops from each bag. Therefore, the number of crops in each bag is intimately tied to both the batch size and the underlying assumptions about the nature (i.e. noise level) of the bags in the data. In most cases, the number of crops in a bag is large, and therefore we must sample mini-bags (e.g. a subset of a bag) to construct a batch. If the bag sizes are too small,

then it is likely that there will not be enough of the true underlying identity in each bag to learn effectively. On the contrary, if the bags are too large, then it is likely that the number of bags in each batch is not sufficient for training. An implicit assumption of this framework is that in expectation, the most common identity in each bag is the identity representative of the bag label. The noise level can still be high without violating this assumption, because most non-representative crops in a bag are of different identities altogether. To ease notation, we will refer to bags and mini-bags interchangeably. In general, we mean mini-bags in algorithmic contexts, and bags in dataset contexts.

During inference, we follow the standard object re-identification procedure. Given a query set, gallery set, and an optional distractor set, we search for a specific object in the gallery based on a query image. All crops are embedded using our instance feature extractor, and then the distance metric used during training is used to return a

ranking over the gallery and distractors for each query image. Based on this ranking, we track the rank-k accuracy for $k \in \{1, 5, 10\}$, as well as the mean average precision (mAP).

---

**Algorithm 1:** Contrastive Multiple Instance Learning (CMIL)

---

**Input:** Set of bags $\mathcal{B} = \{B_1, B_2, ..., B_N\}$, where each bag $B_i$ contains crops $\{x_1^i, x_2^i, ..., x_{M_i}^i\}$

**Output:** Trained model parameters $\theta$ for crop feature extraction

Initialize model parameters $\theta$, $\phi$, and $\psi$ randomly

**while** *not converged* **do**

  **for** $\mathcal{B}_{batch} \subset \mathcal{B}$ **do**

    **for** $B_i \in \mathcal{B}_{batch}$ **do**

      **for** $x_j^i \in B_i$ **do**

        $z_j^i = f_\theta(x_j^i)$ # Extract features from crops

      **end**

      $r_i = g_\phi(\{z_1^i, ..., z_{M_i}^i)\})$ # Aggregate crop features into bag representation

    **end**

    $\mathcal{L}_{triplet}(r_{1,...,|B_{batch}|})$ # Triplet Loss

    $\mathcal{L}_{CE}(h_\psi(r_{1,...,|B_{batch}|})$ # CE Loss

    $\mathcal{L}_{align}(r_{1,...,|B_{batch}|}, z_{1,...,|M_i|}^{1,...,|B_{batch}|})$ # Align Loss

    $\mathcal{L} = \alpha \mathcal{L}_{triplet} + \beta \mathcal{L}_{CE} + \gamma \mathcal{L}_{align}$ # Aggregated Loss

    Update model parameters $\theta$, $\phi$, $\psi$ to minimize $\mathcal{L}$

  **end**

**end**

---

## 3.1 Loss Function

**Table 1: Dataset Summary statistics for each dataset used in the experiments.**

| Dataset | Market-1501 | SYSU-30k | weak MUDD |
|---|---|---|---|
| # identities | 1,501 | 30,508 | 150 |
| Scene | Outdoor | Indoor,Outdoor | Outdoor |
| Annotation | Strong | Weak | Weak |
| Cameras | 6 | Countless | Countless |
| Images | 32,668 | 29,606,918 | 9,069 |

CMIL leverages both the identity and triplet losses. The identity loss is the cross entropy loss when each class represents a person identity

$$\mathcal{L}_{CE} = -\sum_{c=1}^{C} y_c \log(p_c), \tag{4}$$

where $y_c$ is a binary indicator (0 or 1) indicating the label of a sample for class $c$, and $p_c$ is the predicted probability of that class for the same sample, calculated by applying a fully connected layer ($h$ parameterized by $\psi$) and a softmax to the bag representation:

$$p_i = \text{softmax}\Big(h_\psi(r_i)\Big). \tag{5}$$

The triplet loss is:

$$\mathcal{L}_{triplet} = \max\Big(d(r_a, r_p) - d(r_a, r_n) + m_{triplet}, 0\Big) \tag{6}$$

where $r_a$ is a bag representation for an anchor sample, $r_p$ is a bag representation for a positive sample (i.e. a bag with the same label as the anchor sample), and $r_n$ is a bag representation for a negative sample (i.e. a bag with a different label than the anchor sample).

Again, this is explicitly optimizing bag representations and only implicitly optimizing crop representations. In an attempt to address this, we experiment with an *alignment loss*. The intuition is that the most shared identity in a bag is the identity of interest. So an ideal accumulation function is one that can accurately pick out the representative crops, and then create a bag representation very similar to one or all of them. Therefore, we create the alignment loss to encourage the bag representation for a bag $a$ with crops $\{x_1^a, x_2^a, \ldots, x_{N_a}^a\}$ to be close to any one of the crop representations:

$$\mathcal{L}_{align} = \max\Big(0, \min\{d(r_a, z_a^1), d(r_a, z_a^2), \ldots, d(r_a, z_a^{N_i})\} - m_{align}\Big), \tag{7}$$

where $m_{align}$ is a margin hyperparameter.

The total loss function is a weighted combination of the identity, triplet, and alignment losses. The weighting for each loss (i.e. $\alpha$, $\beta$, and $\gamma$) is selected during our hyperparameter search.

$$\mathcal{L} = \alpha \mathcal{L}_{triplet} + \beta \mathcal{L}_{CE} + \gamma \mathcal{L}_{align} \tag{8}$$

Note that the triplet loss can be substituted with any contrastive loss. Algorithm 1 provides an overview of CMIL in pseudocode.

## 4 EXPERIMENTS

We evaluate our methodology on three datasets:

- **WL-Market-1501:** The widely used Market-1501 person ReID dataset [55], but with synthetically weak labels. The synthetic labels are generated by duplicating images from the training set some number of times, and assigning them to random bags.
- **WL-MUDD:** Our real-world dataset introduced in Section 2.2
- **SYSU30k:** A large-scale weakly supervised person ReID dataset with over 29 million images gathered from TV program videos. The videos are randomly broken into clips, and then each clip is manually annotated with an identity, but all detected people are noisily assigned that identity, forming bag-level labels.

The dataset statistics can be seen in Table 1. While the training set of each of these datasets is weakly labeled, the test sets are accurately labeled for normal person ReID evaluation [50]. We track the mean average precision (mAP) and the Rank-k accuracy for $k \in \{1, 5, 10\}$.

## 4.1 Implementation Details and Hyperparameter Tuning

## 4.2 Baseline Methods

Ye et al. [49] introduce online CO-REfining (CORE), a framework for online co-refining of ReID models. CORE uses learning rate schedules to optimize two models collaboratively, while also iteratively refining the noisy labels in a dataset. CORE is a state-of-the-art

**Table 2: The sweep configuration for hyperparameter optimization, along with the final CMIL hyperparameters for each dataset. $U_{int}(x, y)$ represents an integer uniform distribution from $x$ to $y$, $U_{log}(x, y)$ represents a log uniform, and $U(x, y)$ represents a standard uniform distribution on all real numbers from $x$ to $y$.**

| Parameter | Search Range | Final Values | | |
|---|---|---|---|---|
| | | Market-1501 | SYSU-30k | Weak MUDD |
| bag size | $U_{int}(5, 10)$ | 6 | 5 | 9 |
| batch size | $U_{int}(5, 10)$ | 10 | 10 | 5 |
| distance metric | [euclidean, cosine] | cosine | cosine | cosine |
| fixbase epoch | $U_{int}(0, 10)$ | 7 | 10 | 8 |
| learning rate | $U_{log}(1e-05, 0.01)$ | 2.1153e-4 | 2.828e-3 | 4.044e-4 |
| margin | $U(0.1, 1)$ | 0.9992 | 0.8592 | 0.7731 |
| feature norm | [false, true] | False | False | False |
| gamma | [0, 0.01, 0.1] | 0 | 0 | 0 |
| alpha | $U(0, 1)$ | 0.5638 | 0.8083 | 0.3882 |
| beta | $U(0, 1)$ | 0.3872 | 0.9242 | 0.7339 |

**Table 3: Results on the WL-Market1501 dataset at varying levels of noise. The noise level represents the percentage of the dataset with incorrect labels. This dataset was synthetically constructed by duplicating images in the training set and assigning them to random bags – 75% noise would correspond to duplicating each image three times, therefore only 1 in 4 images would be correctly labeled.**

| Noise | Method | R1 | R5 | R10 | mAP |
|---|---|---|---|---|---|
| | CORE | 80.9% | 92.2% | 95.0% | 48.6% |
| 50% | MIML | 71.8% | 87.0% | 91.5% | 46.3% |
| | CMIL (Ours) | 80.7% | 91.9% | 94.4% | 56.8% |
| | CORE | 68.1% | 83.6% | 88.2% | 38.6% |
| 66% | MIML | 62.7% | 82.2% | 87.6% | 38.8% |
| | CMIL (Ours) | 76.4% | 89.6% | 93.0% | 54.4% |
| | CORE | 56.1% | 74.4% | 80.8% | 27.9% |
| 75% | MIML | 50.4% | 71.6% | 79.31% | 26.6% |
| | CMIL (Ours) | 70.0% | 86.4% | 90.9% | 48.8% |
| | CORE | 47.5% | 66.0% | 73.2% | 17.4% |
| 80% | MIML | 54.0% | 60.8% | 71.1% | 19.0% |
| | CMIL (Ours) | 64.9% | 82.8% | 88.0% | 43.9% |

method for learning ReID models among noisy labels and weak supervision.

Meng et al. [29] introduce Cross View Multi-Instance Multi-Label Learning (CV-MIML). Being based on MIL, this method falls most closely related to ours. Although originally developed for the setting where a target person is known to appear within an untrimmed video but no further information is available, this weakly supervised setting is equivalent to ours, although perhaps simpler due to correlations within a single video frame. Importantly, this method only performs bag classification during training, taking advantage only of the identity loss. Instead, CMIL optimizes bag representations explicitly.

We implement the CMIL framework using PyTorch. For a fair comparison, all methods utilize ResNet-50, pretrained on Imagenet [6],

as the feature extractor $f_\theta$. Our method also requires a reduction function $g$, and in this case, we use a 2-layer set transformer [20]. Section 5.1 includes ablations where we experiment with simpler reduction functions, namely the average, max, and sum operators. Importantly, we also implement a *bag* sampling function. We expect two conditions to be met for every mini-batch:

(1) Each batch will consist of $b$ sub-bags, where a sub-bag is a subset of a bag. If a bag is smaller than $b$, then the bag is oversampled.
(2) Each bag label present in the mini-batch will have two or more bags in the mini-batch to ensure that valid triplets can always be constructed.

For hyperparameter selection, we run a Bayes hyperparameter search with early stopping (if model validation accuracy has not improved in 5 epochs, terminate the run) and hyperband (with an eta value of 2 and a minimum iteration count of 3) for early termination of less promising runs [21]. Table 2 describes the hyperparameter search ranges. The search aims to maximize the rank-1 accuracy on the validation set over 50 epochs. For each dataset, 250 models with hyperparameters sampled from the listed distributions were trained and evaluated, and the best-performing hyperparameters are also shown in Table 2. Finally, using the best-performing hyperparameters, a final training run was done using the combined training and validation set, evaluated on the test set, and reported in our results.

## 5 RESULTS AND DISCUSSION

Table 3 summarizes the rank-1 (R1), rank-5 (R5), rank-10 (R10) accuracy and mean average precision (mAP) of the different methods on the Market-1501 dataset with varying levels of synthetic label noise. At the 50% noise level, our CMIL method achieves an R1 accuracy of 80.7%, nearly matching the performance of CORE and outperforming MIML by 8.9%. As the noise level begins to increase, CMIL dominates the other methods by a growing margin. At 80% label noise, the hardest setting, CMIL obtains a R1 of 64.9%, which signifies a 10.9% boost over the best baseline. The consistent gaps between CMIL and other approaches illustrate that our method

can effectively learn useful representations among even extremely noisy bags.

Table 4 summarizes the performance of different methods when trained on the real-world Weak MUDD dataset. With noisy group annotations, our CMIL framework obtains 73.2% rank-1 accuracy. This significantly outperforms baseline methods, including CORE and MIML, by 2.5% and 6% respectively.

**Table 4: Results on the WL-MUDD dataset.**

| Method | R1 | R5 | R10 | mAP |
|--------|------|------|------|------|
| CORE | 67.2% | 83.3% | 92.3% | 71.6% |
| MIML | 70.7% | 87.7% | 95.2% | 74.6% |
| CMIL (Ours) | 73.2% | 90.0% | 96.8% | 75.1% |

**Table 5: Results on the SYSU30k dataset.**

| Supervision | Method | R1 |
|-------------|--------|------|
| Transfer Learning | DARIR [38] | 11.2% |
| | DF [7] | 10.3% |
| | Local CNN [46] | 23.0% |
| | MGN [41] | 23.6% |
| Self-Supervised | SimCLR [2] | 10.9% |
| | MoCo v2 [4] | 11.6% |
| | BYOL [13] | 12.7% |
| | Triplet [40] | 27.5% |
| Weakly Supervised | W-Local CNN [39] | 28.8% |
| | W-MGN [39] | 29.5% |
| | WS-TAL [26] | 34.4% |
| | CMIL (Ours) | 33.9% |

The SYSU-30k dataset is very large and computationally expensive to optimize models on. Therefore, we compare directly to the results reported in prior work in Table 5. CMIL attains 33.9% R1 accuracy outperforming the best transfer and self-supervised learning approaches by 10.3% and 6.4% respectively. The best weakly supervised method is WS-TAL [26], which is specifically engineered to optimize ReID models when the labels are generated from video tracklets, matching the SYSU construction. WS-TAL reaches 34.4% R1 accuracy. CMIL nearly matches this performance, lagging by only 0.5%, using more general labeling assumptions.

Surprisingly, in every case, the alignment loss does not improve accuracy – as shown in Table 2, $\gamma = 0$ and therefore the alignment loss is not used. During training, CMIL models are optimized at the bag level. Given a batch of bag representations, the model is optimized for bags with the same label to be close, and bags with a different label to be far from each other in representation space. The bag representations are built from crop representations, but nothing is preventing the bag and crop representations from being far apart. This should be problematic, because at test time, we are evaluating the quality of the crop embeddings.

Figure 4 plots the rank-1 accuracy and alignment loss versus training step when $\gamma = 0$ (i.e. the alignment loss is not used). We
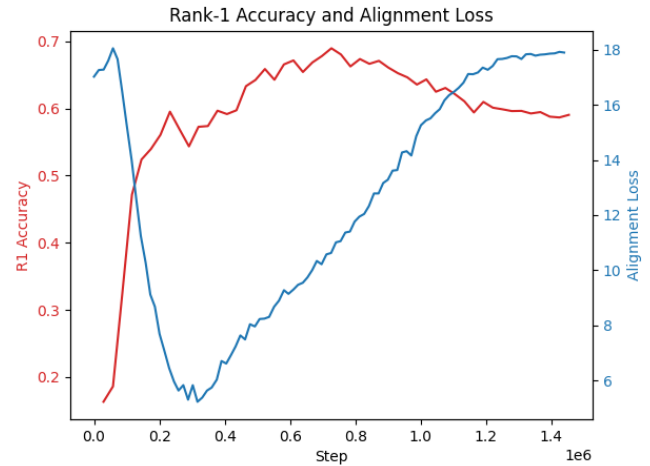


**Figure 4: The rank-1 accuracy and the alignment loss throughout a training run. The alignment loss exhibits unintuitive behavior - the best alignment (i.e. lowest) does not correspond to the best model accuracy (i.e. highest). This behavior is characteristic of every model trained in this work, including those using different accumulation functions.**

see that the bag and crop representations start quite different, and then begin growing more similar. However, there reaches a point relatively early in training where the alignment between instance and bag representations begins to decrease. Interestingly, the performance of the model (therefore the *crop* embedding model) continues to improve, even while their embeddings diverge from bag embeddings. This phenomenon is consistent across all of our experiments.

It is not likely that this counterintuitive behavior is an artifact of the way we are measuring the alignment loss. It is not known which crop within a bag is the crop representative of the bag label, so the alignment loss used here is the distance, using the same measure used during training, between the bag representation and the *closest* crop representation within the bag. A better understanding of this phenomenon requires further research.

## 5.1 Ablation Study

In this ablation, we experiment with other, more traditional, choices for permutation invariant accumulation function, specifically the max, average, and sum. Each bag contains crops, and one or more of the crops in the bag are representative of the bag label. Of course, which crop specifically is unknown. Intuitively, the job of the accumulation function is to select the crop (or a representation of the collection of crops) that corresponds to the bag label.

All aforementioned models have used a set transformer for the accumulation function, as the learnable attention-based model makes it possible to behave as a selector, or any arbitrary combination of the crop representations. However, it does come at the cost of complexity. Other reasonable, and much simpler, choices are to set the bag representation to be the max, average, or sum of the crop representations.

**Table 6: Comparison of different accumulation functions on WL-Market-1501 and WL-MUDD datasets. Using a simple average of crop representations performs nearly as well as the set transformer.**

|  | WL-Market1501 | | | | WL-MUDD | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| Method | R1 | R5 | R10 | mAP | R1 | R5 | R10 | mAP |
| CMIL w/Set Transformer | 70.0% | 86.4% | 90.9% | 48.8% | 73.2% | 90.0% | 96.8% | 75.1% |
| CMIL w/Max | 60.2% | 78.7% | 84.6% | 33.9% | 66.8% | 82.1% | 90.7% | 68.2% |
| CMIL w/Avg | 69.8% | 85.6% | 90.1% | 44.1% | 71.1% | 88.6% | 94.8% | 74.5% |
| CMIL w/Sum | 51.2% | 72.0% | 79.6% | 24.3% | 64.3% | 79.3% | 88.9% | 62.4% |

In Table 6, we compare the performance of the different accumulation functions on both the WL-Market1501 and the WL-MUDD datasets. Interestingly, the average does well, matching, or nearly matching, the performance of the set transformer. This is surprising as the crops within a representative bag of the bag label are the minority of samples, typically representing less than half of the samples within each bag. This could indicate that the set transformer is roughly just performing an average. A potential reason for this is that the non-representative crop features could act to cancel one another out such that the bag representation is still close to the corresponding representative features. Interestingly, tracking the alignment loss for each of these simpler accumulation functions shows the exact behavior as depicted in Figure 4.

## 6   RELATED WORK

A large body of work has focused on supervised re-ID, where models are trained on data with individual object identity labels [17, 23, 50, 56, 57]. These approaches employ deep neural networks, to extract visual features that are representative of specific identities, and optimize them according to various loss functions, including the identification loss where each identity is treated as a class [57], verification loss where pairwise relationships are optimized vi the contrastive loss [34], triplet loss that treats the problem as a retrieval ranking problem [18], and others have been proposed to optimize re-ID performance [1, 3, 15, 27, 31, 36, 45, 47, 59]. These methods perform well on baseline datasets, where ample data labeling is available.

To alleviate the labeling bottleneck, recent works have begun investigating re-ID under weak supervision. Strategies include exploiting image-level labels [29], pseudo-labels [39], noisy label refinement [49], online captions [14, 54], and domain adaptation [52]. While showing promise, these methods still fall behind those that are fully supervised [58].

Most similarly to our work, Meng et al. [29] leverage image-level labels in conjunction with Multiple Instance Learning (MIL) to effective facial recognition models. MIL offers a paradigm to handle label ambiguity in training data by modeling labels at a bag level. A bag can be a collection of instances associated with a particular label, but we only know that one or more of the instances in that collection truly belongs to that label. Several works have adapted this specifically for treating video "tracklets" as a bag of instances [26, 42] MIL has found diverse applications including image classification (particularly medical imagery) [32, 43], object detection [19, 35, 53], and drug discovery [11]. Critically, Meng et al. [29] apply MIL to the person re-identification problem in the

identity loss setting. In contrast, CMIL improves upon this by allowing for use of contrastive learning, which has shown significant advantages in person ReID and related settings [12, 18].

Lastly, we must mention the work in unsupervised ReID [8, 9, 22, 24, 37, 51]. These methods do not require labels. Typically, these methods use iterative clustering and classification, such that unlabeled images are clustered into "pseudo" classes, which are then used to train or update a model. Then the new/updated model is used to refine the pseudo labels, and so on. Improvements to this standard approach include substituting the clustering step for pairwise comparisions [25], and an improved clustering step by improving the global clusters using ensembles of image-part based predictions [5]. Of course, performance is still greatly improved when labels are present [10, 28, 44, 48, 60].

## 7   CONCLUSION

In this paper, we introduced Contrastive Multiple Instance Learning (CMIL), a novel framework tailored for more effective person re-identification under weak supervision. CMIL tackles the challenge of learning discriminative person representations when only bag-level labels indicating a shared identity among a group of photos are available. Although the model is trained at the bag level, the person-level representations improve alongside the quality of the bag-level representations. We experiment with adding an alignment loss to further encourage the person and bag representations to be similar, but found it ineffective empirically.

We experiment on three datasets, one of which is the Weakly Labeled Muddy Racer Re-Identification Dataset (WL-MUDD), which is curated and released from real-world weak labels from PerformancePhoto.co. Across these experiments, CMIL consistently achieved state-of-the-art rank-1, rank-5, and rank-10 accuracy as well as mean average precision. On the large-scale SYSU-30k dataset, CMIL matched the top-reported result while requiring fewer assumptions. Ablations also revealed surprising effectiveness of average pooling for instance aggregation, suffering only slight performance degradation to the set transformer.

The contributions of this work are threefold. First, we introduce the new Weakly Labeled Muddy Racer Re-Identification dataset (WL-MUDD) built from PerformancePhoto.co, an off-road photograph platform. Second, we introduce the CMIL framework that enables efficient exploitation of cheap weak supervision for person re-id through enabling contrastive learning with Multiple Instance Learning. And third, we show the efficacy of CMIL on two real-world datasets and one synthetic, outperforming baselines.

# REFERENCES

[1] Guangyi Chen, Chunze Lin, Liangliang Ren, Jiwen Lu, and Jie Zhou. 2019. Self-critical attention learning for person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 9637–9646.

[2] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. 2020. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*. PMLR, 1597–1607.

[3] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. 2017. Beyond triplet loss: a deep quadruplet network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 403–412.

[4] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. 2020. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297* (2020).

[5] Yoonki Cho, Woo Jae Kim, Seunghoon Hong, and Sung-Eui Yoon. 2022. Part-based pseudo label refinement for unsupervised person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7308–7318.

[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.

[7] Shengyong Ding, Liang Lin, Guangrun Wang, and Hongyang Chao. 2015. Deep feature learning with relative distance comparison for person re-identification. *Pattern Recognition* 48, 10 (2015), 2993–3003.

[8] Hehe Fan, Liang Zheng, Chenggang Yan, and Yi Yang. 2018. Unsupervised person re-identification: Clustering and fine-tuning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 14, 4 (2018), 1–18.

[9] Dengpan Fu, Dongdong Chen, Jianmin Bao, Hao Yang, Lu Yuan, Lei Zhang, Houqiang Li, and Dong Chen. 2021. Unsupervised Pre-Training for Person Re-Identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 14750–14759.

[10] Dengpan Fu, Dongdong Chen, Hao Yang, Jianmin Bao, Lu Yuan, Lei Zhang, Houqiang Li, Fang Wen, and Dong Chen. 2022. Large-Scale Pre-Training for Person Re-Identification With Noisy Labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2476–2486.

[11] Gang Fu, Xiaofei Nan, Haining Liu, Ronak Y Patel, Pankaj R Daga, Yixin Chen, Dawn E Wilkins, and Robert J Doerksen. 2012. Implementation of multiple-instance learning in drug activity prediction. In *BMC bioinformatics*, Vol. 13. BioMed Central, 1–12.

[12] Saurabh Garg, Amrith Setlur, Zachary Chase Lipton, Sivaraman Balakrishnan, Virginia Smith, and Aditi Raghunathan. 2023. Complementary Benefits of Contrastive Learning and Self-Training Under Distribution Shift. *arXiv preprint arXiv:2312.03318* (2023).

[13] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. 2020. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems* 33 (2020), 21271–21284.

[14] Matthieu Guillaumin, Jakob Verbeek, and Cordelia Schmid. 2010. Multiple instance metric learning from automatically labeled bags of faces. In *Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, Greece, September 5-11, 2010, Proceedings, Part I 11*. Springer, 634–647.

[15] Jianyuan Guo, Yuhui Yuan, Lang Huang, Chao Zhang, Jin-Ge Yao, and Kai Han. 2019. Beyond human parts: Dual part-aligned representations for person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 3642–3651.

[16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[17] Shuting He, Hao Luo, Pichao Wang, Fan Wang, Hao Li, and Wei Jiang. 2021. Transreid: Transformer-based object re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 15013–15022.

[18] Alexander Hermans, Lucas Beyer, and Bastian Leibe. 2017. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737* (2017).

[19] Fang Huang, Jinqing Qi, Huchuan Lu, Lihe Zhang, and Xiang Ruan. 2017. Salient object detection via multiple instance learning. *IEEE Transactions on Image Processing* 26, 4 (2017), 1911–1922.

[20] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. 2019. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International conference on machine learning*. PMLR, 3744–3753.

[21] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. 2018. Hyperband: A novel bandit-based approach to hyperparameter optimization. *Journal of Machine Learning Research* 18, 185 (2018), 1–52.

[22] Minxian Li, Xiatian Zhu, and Shaogang Gong. 2018. Unsupervised Person Re-identification by Deep Learning Tracklet Association. In *Proceedings of the European Conference on Computer Vision (ECCV)*.

[23] Wei Li, Xiatian Zhu, and Shaogang Gong. 2018. Harmonious attention network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2285–2294.

[24] Yutian Lin, Xuanyi Dong, Liang Zheng, Yan Yan, and Yi Yang. 2019. A bottom-up clustering approach to unsupervised person re-identification. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 33. 8738–8745.

[25] Yutian Lin, Lingxi Xie, Yu Wu, Chenggang Yan, and Qi Tian. 2020. Unsupervised person re-identification via softened similarity learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3390–3399.

[26] Min Liu, Yuan Bian, Qing Liu, Xueping Wang, and Yaonan Wang. 2023. Weakly Supervised Tracklet Association Learning with Video Labels for Person Re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).

[27] Zimo Liu, Jingya Wang, Shaogang Gong, Huchuan Lu, and Dacheng Tao. 2019. Deep reinforcement active learning for human-in-the-loop person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 6122–6131.

[28] Hao Luo, Pichao Wang, Yi Xu, Feng Ding, Yanxin Zhou, Fan Wang, Hao Li, and Rong Jin. 2021. Self-supervised pre-training for transformer-based person re-identification. *arXiv preprint arXiv:2111.12084* (2021).

[29] Jingke Meng, Sheng Wu, and Wei-Shi Zheng. 2019. Weakly supervised person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 760–769.

[30] Jingke Meng, Wei-Shi Zheng, Jian-Huang Lai, and Liang Wang. 2021. Deep graph metric learning for weakly supervised person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 10 (2021), 6074–6093.

[31] Jifei Song, Yongxin Yang, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. 2019. Generalizable person re-identification by domain-invariant mapping network. In *Proceedings of the IEEE/CVF conference on Computer Vision and Pattern Recognition*. 719–728.

[32] PJ Sudharshan, Caroline Petitjean, Fabio Spanhol, Luiz Eduardo Oliveira, Laurent Heutte, and Paul Honeine. 2019. Multiple instance learning for histopathological breast cancer image classification. *Expert Systems with Applications* 117 (2019), 103–111.

[33] Jacob Tyo, Motolani Olarinre, Youngseog Chung, and Zachary C Lipton. 2023. MUDD: A New Re-Identification Dataset with Efficient Annotation for Off-Road Racers in Extreme Conditions. *arXiv preprint arXiv:2311.08488* (2023).

[34] Rahul Rama Varior, Bing Shuai, Jiwen Lu, Dong Xu, and Gang Wang. 2016. A siamese long short-term memory architecture for human re-identification. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part VII 14*. Springer, 135–153.

[35] Fang Wan, Chang Liu, Wei Ke, Xiangyang Ji, Jianbin Jiao, and Qixiang Ye. 2019. C-mil: Continuation multiple instance learning for weakly supervised object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2199–2208.

[36] Cheng Wang, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. 2018. Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In *Proceedings of the European conference on computer vision (ECCV)*. 365–381.

[37] Dongkai Wang and Shiliang Zhang. 2020. Unsupervised Person Re-Identification via Multi-Label Classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[38] Guangrun Wang, Liang Lin, Shengyong Ding, Ya Li, and Qing Wang. 2016. DARI: Distance Metric and Representation Integration for Person Verification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 30.

[39] Guangrun Wang, Guangcong Wang, Xujie Zhang, Jianhuang Lai, Zhengtao Yu, and Liang Lin. 2020. Weakly supervised person re-id: Differentiable graphical learning and a new benchmark. *IEEE Transactions on Neural Networks and Learning Systems* 32, 5 (2020), 2142–2156.

[40] Guangrun Wang, Keze Wang, Guangcong Wang, Philip HS Torr, and Liang Lin. 2021. Solving inefficiency of self-supervised representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9505–9515.

[41] Guanshuo Wang, Yufeng Yuan, Xiong Chen, Jiwei Li, and Xi Zhou. 2018. Learning discriminative features with multiple granularities for person re-identification. In *Proceedings of the 26th ACM international conference on Multimedia*. 274–282.

[42] Xueping Wang, Min Liu, Dripta S Raychaudhuri, Sujoy Paul, Yaonan Wang, and Amit K Roy-Chowdhury. 2021. Learning person re-identification models from videos with weak supervision. *IEEE Transactions on Image Processing* 30 (2021), 3017–3028.

[43] Jiajun Wu, Yinan Yu, Chang Huang, and Kai Yu. 2015. Deep multiple instance learning for image classification and auto-annotation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3460–3469.

[44] Suncheng Xiang, Dahong Qian, Jingsheng Gao, Zirui Zhang, Ting Liu, and Yuzhuo Fu. 2023. Rethinking person re-identification via semantic-based pretraining. *ACM Transactions on Multimedia Computing, Communications and Applications* 20, 3 (2023), 1–17.

[45] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. 2017. Joint detection and identification feature learning for person search. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3415–3424.

[46] Jiwei Yang, Xu Shen, Xinmei Tian, Houqiang Li, Jianqiang Huang, and Xian-Sheng Hua. 2018. Local convolutional neural networks for person re-identification. In *Proceedings of the 26th ACM international conference on Multimedia*. 1074–1082.

[47] Qize Yang, Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. 2019. Patch-based discriminative feature learning for unsupervised person re-identification. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3633–3642.

[48] Zizheng Yang, Xin Jin, Kecheng Zheng, and Feng Zhao. 2022. Unleashing potential of unsupervised pre-training with intra-identity regularization for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14298–14307.

[49] Mang Ye, He Li, Bo Du, Jianbing Shen, Ling Shao, and Steven CH Hoi. 2021. Collaborative refining for person re-identification with label noise. *IEEE Transactions on Image Processing* 31 (2021), 379–391.

[50] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. 2021. Deep learning for person re-identification: A survey and outlook. *IEEE transactions on pattern analysis and machine intelligence* 44, 6 (2021), 2872–2893.

[51] Hong-Xing Yu, Wei-Shi Zheng, Ancong Wu, Xiaowei Guo, Shaogang Gong, and Jian-Huang Lai. 2019. Unsupervised Person Re-Identification by Soft Multilabel Learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[52] Jongmin Yu, Hyeontaek Oh, Minkyung Kim, and Junsik Kim. 2023. Weakly supervised contrastive learning for unsupervised vehicle reidentification. *IEEE transactions on neural networks and learning systems* (2023).

[53] Tianning Yuan, Fang Wan, Mengying Fu, Jianzhuang Liu, Songcen Xu, Xiangyang Ji, and Qixiang Ye. 2021. Multiple instance active learning for object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5330–5339.

[54] Shizhen Zhao, Changxin Gao, Yuanjie Shao, Wei-Shi Zheng, and Nong Sang. 2021. Weakly supervised text-based person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 11395–11404.

[55] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. 2015. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*. 1116–1124.

[56] Liang Zheng, Yi Yang, and Alexander G Hauptmann. 2016. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984* (2016).

[57] Liang Zheng, Hengheng Zhang, Shaoyan Sun, Manmohan Chandraker, Yi Yang, and Qi Tian. 2017. Person re-identification in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1367–1376.

[58] Mingkai Zheng, Fei Wang, Shan You, Chen Qian, Changshui Zhang, Xiaogang Wang, and Chang Xu. 2021. Weakly supervised contrastive learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 10042–10051.

[59] Sanping Zhou, Fei Wang, Zeyi Huang, and Jinjun Wang. 2019. Discriminative feature learning with consistent attention regularization for person re-identification. In *Proceedings of the IEEE/CVF international conference on computer vision*. 8040–8049.

[60] Kuan Zhu, Haiyun Guo, Tianyi Yan, Yousong Zhu, Jinqiao Wang, and Ming Tang. 2022. Pass: Part-aware self-supervised pre-training for person re-identification. In *European Conference on Computer Vision*. Springer, 198–214.