

Tasks Integrated Networks: Joint Detection and Retrieval for Image Search

Lei Zhang, Zhenwei He, Yi Yang, Liang Wang, Xinbo Gao

Abstract—The traditional object (person) retrieval (re-identification) task aims to learn a discriminative feature representation with intra-similarity and inter-dissimilarity, which supposes that the objects in an image are manually or automatically pre-cropped exactly. However, in many real-world searching scenarios (e.g., video surveillance), the objects (e.g., persons, vehicles, etc.) are seldom accurately detected or annotated. Therefore, object-level retrieval becomes intractable without bounding-box annotation, which leads to a new but challenging topic, i.e. image-level search with multi-task integration of joint detection and retrieval. In this paper, to address the image search issue, we first introduce an end-to-end Integrated Net (I-Net), which has three merits: 1) A Siamese architecture and an on-line pairing strategy for similar and dissimilar objects in the given images are designed. Benefited by the Siamese structure, I-Net learns the shared feature representation, because, on which, both object detection and classification tasks are handled. 2) A novel on-line pairing (OLP) loss is introduced with a dynamic feature dictionary, which alleviates the multi-task training stagnation problem, by automatically generating a number of negative pairs to restrict the positives. 3) A hard example priority (HEP) based softmax loss is proposed to improve the robustness of classification task by selecting hard categories. The shared feature representation of I-Net may restrict the task-specific flexibility and learning capability between detection and retrieval tasks. Therefore, with the philosophy of **divide and conquer**, we further propose an improved I-Net, called DC-I-Net, which makes two new contributions: 1) two modules are tailored to handle different tasks separately in the integrated framework, such that the task specification is guaranteed. 2) A class-center guided HEP loss (C²HEP) by exploiting the stored class centers is proposed, such that the intra-similarity and inter-dissimilarity can be captured for ultimate retrieval. Extensive experiments on famous image-level search oriented benchmark datasets, such as CUHK-SYSU dataset and PRW dataset for person search and the large-scale WebTattoo dataset for tattoo search, demonstrate that the proposed DC-I-Net outperforms the state-of-the-art tasks-integrated and tasks-separated image search models.

Index Terms—Image Search, Object Detection, Re-identification, Retrieval, Deep Learning

1 INTRODUCTION

SEARCHING images containing some interested object from a large gallery image set is a new but challenging research area. For example, in real-world video surveillance, many tasks such as criminals search [1] and multi-camera tracking [2] are closely related to our life. These tasks need to search the image containing a target person from the videos of different scenes and backgrounds where the persons in the videos are not cropped or annotated. Generally, for person search tasks, the machine should first know where are the persons in the image (pedestrian detection) and then guess who is the right person (person re-identification). Therefore, image search problem is closely related to two independent computer vision tasks, such as object detection (positioning) and object retrieval (matching). The detection model aims to locate the interested objects in the images, while the purpose of the retrieval model is to match the query objects

and gallery objects where the two images may come from different distributions. For example, in pedestrian detection and person re-identification problems, camera views, poses, occlusions, illuminations, backgrounds and resolutions may easily cause intra-class dissimilarity and inter-class similarity. Therefore, both detection and retrieval are challenging problems in computer vision and have attracted lots of attention in recent years [3], [4], [5], [6], [7], [8].

Person search is the first attempt for image search issue. Before that, although numerous endeavor on person detection and re-identification has been made, most of them handle these two problems independently. That is, the traditional methods divide the person search task into two separated subtasks. First, a pedestrian detector is implemented to predict the bounding boxes of persons from images, and then the Rectangular regions of detected persons are cropped based on the coordinates of the predicted bounding boxes. Second, the feature representations of the detected person regions of interest are computed for person re-identification. In general person re-identification (Re-ID) task, the pedestrian images are manually annotated and cropped for discriminative feature representation network training [3], [9], [10], which gives rise to two considerations. On one hand, in real-world video surveillance task, most of the detectors inevitably have the false alarms and misalignments, which, to some extent, may cause a significant performance drop of re-identification accuracy. On the other hand, the two independent subtasks seem to be less user-friendly for ultimate Re-ID in real applications. Therefore,

- L. Zhang and Z. He are with the School of Microelectronics and Communication Engineering, Chongqing University, Chongqing 400044, China. (E-mail: leizhang@cqu.edu.cn, hzw@cqu.edu.cn).
- Y. Yang is with the Center for Artificial Intelligence, University of Technology Sydney, Australia. (E-mail: yi.yang@uts.edu.au).
- L. Wang is with the National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China. (E-mail: wangliang@nlpr.ia.ac.cn).
- X. Gao is with the Chongqing Key Laboratory of Image Cognition, Chongqing University of Posts and Telecommunications, Chongqing 400065. (E-mail: gaoyb@cqupt.edu.cn).

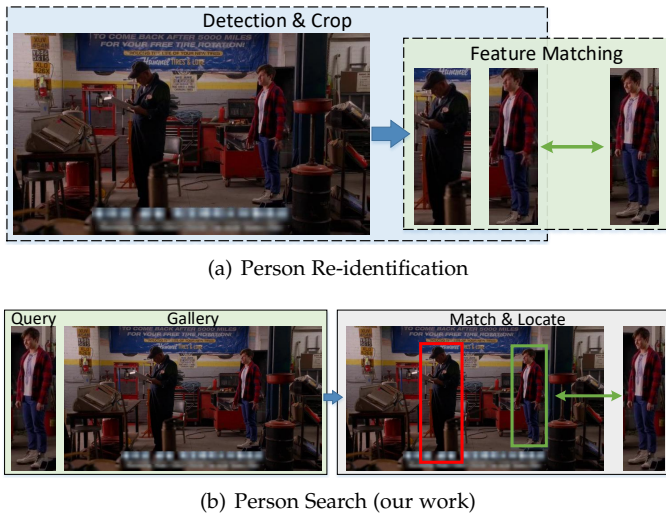


Fig. 1. Comparison of person re-identification and person search. In (a), the traditional person re-identification task needs to first detect and crop all person regions for feature matching. The person search task in (b) means the joint detection and re-identification for localization.

in order to reduce the gap between the traditional person Re-ID and real-world application, we propose to deal with a person search task by jointly detecting and matching the target person when given only two images, which is more user-friendly. Person search task focuses on the co-learning of detection and person re-identification, which means the two tasks can adapt to each other in an integrated framework to outperform single-task models. The difference between traditional person re-identification and person search is clearly presented in Fig. 1.

Specifically, we propose an Integrated Net (I-Net) which learns the detector and re-identifier in a deep unified framework end-to-end for image search task. Three important aspects should be taken into consideration. (1) a detector should be tailored for locating the objects given an image. In I-Net, the Faster-RCNN [11] based two-stage detection mechanism is considered. (2) a feature representation should be learned for object retrieval. However, the traditional metric learning loss such as triplet loss [12] or triplet-wise architecture are hard to be directly implemented on the end-to-end detection and retrieval integration training structure due to the lack of objects (e.g., persons) with different identities in each iteration caused by the inherent few input images of detection task. To this end, a softmax guided online pairing (OLP) metric loss and hard example priority (HEP) classification loss are proposed for intra-similarity and inter-dissimilarity learning of feature representation in a pair-wise Siamese architecture. The joint learning of the softmax guided metric loss (OLP) and classification (HEP) loss helps to learn more discriminative representations that benefit to the person search task. (3) the detector and re-identifier should be jointly trained, because the co-learning of two modules can promote the adaptation capability of each other, the retrieval performance and simultaneously the user-friendly property in real-world application.

The structure of our I-Net is based on the Siamese network, as is shown in Fig. 2. The two-stream network

structure gives rise to three advantages:

1) With a Siamese structure, a pair of images can be fed into the model, such that effective feature representation with intra-similarity and inter-dissimilarity for object retrieval can be trained via joint metric and classification loss functions. More importantly, the number of input images can be enlarged, which is very helpful for better learning of the detection module (i.e., Faster-RCNN [11]). Otherwise, the retrieval module will be difficult to be trained without successful detection. Additionally, with more input images, more positive pairs can be constructed to improve the metric learning of the retrieval module with better features.

2) The OLP loss with a dynamic and on-line feature dictionary is proposed, in which the stored features can help the loss function to generate negative pairs to restrict the positive pairs. Comparing to the famous triplet loss [12] for similarity metric learning, which is widely used in the traditional person re-identification tasks [13], [14], [15], [16], our loss can well remit the stagnation problem that is easily encountered in training the triplet loss. The stagnation is due to that for the end-to-end person search task, the input image number of each iteration is so small that the number of identities for generating the positive and negative pairs is not enough for training. In our OLP loss, an on-line feature dictionary is deployed to solve the scarcity of samples for training. Additionally, the OLP is designed with a cross-entropy formulation for similarity metric learning by computing the confidence probability of similarity. The stagnation problem can be well remitted.

3) The HEP loss is further introduced as a classification loss, which is complementary to the metric loss, i.e., OLP, because HEP loss improves the class (identity) discrimination ability of the learned network. Different from the traditional classification loss, the HEP loss selects and focuses on the hard categories via a hard example priority strategy. With the integration of the newly designed OLP and HEP for retrieval and the standard detector losses (i.e., cls. vs. reg.) in Faster-RCNN [11] for detection, the proposed I-Net can be trained end-to-end for both detection and retrieval tasks, and user-friendly image search is achieved.

New Challenges and Improved Solutions of I-Net. As shown in Fig. 2, the shared feature representation of the FC-layer is learned between detection and retrieval in I-Net, which is a little defective due to that the features for detection are focused on the discrimination between foreground and background, while the features for retrieval are focused on the discrimination among all the foregrounds. Therefore, utilizing the shared feature representation from the same layer for both detection and retrieval tasks may deteriorate the performance of image search. Also, in I-Net, the object proposals rather than precise objects are used for retrieval, which may also degrade the performance. With the above new challenges, in this paper, we revisit the image search framework based on I-Net, and further propose new solutions by extending our I-Net [17] with reconsideration of the feature shared network structure problem for joint multi-task co-training. Specifically, with the philosophy of divide and conquer, in the I-Net, we naturally deploy different layered features for each task rather than the shared representation. Specifically, an improved divide and conquer I-Net with a new network structure and improved loss functions,

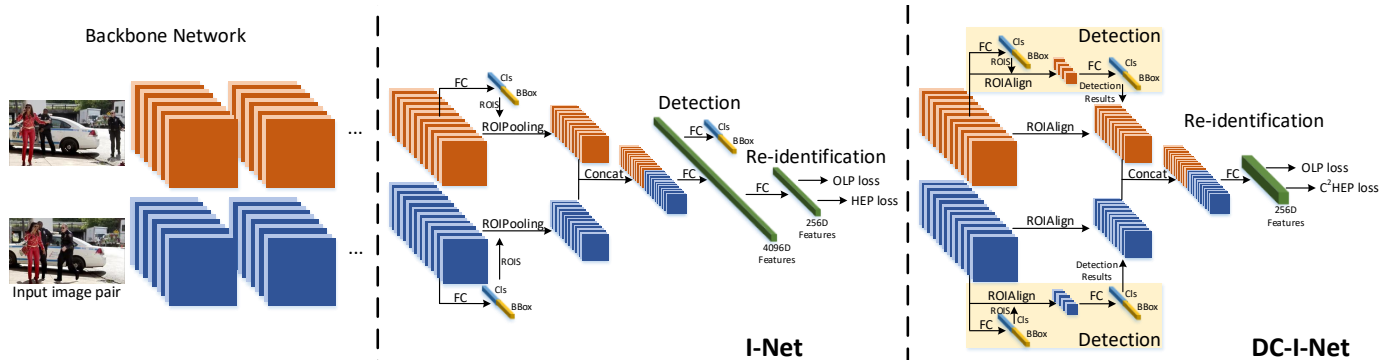


Fig. 2. The flowchart of the proposed models. The Siamese network structure of the proposed I-Net (middle) and DC-I-Net (right) are presented, for both of which, the same backbone network (left) with shared weights is used for feature maps extraction. A pair of images containing the same objects (e.g., persons) are fed into the model for training. For I-Net, two Region Proposal Networks (i.e., Pedestrian Proposal Networks in this case) are implemented for getting the object proposals in each image, and the proposal features generated from the ROI-pooling layers are concatenated and fed into the fully-connected layers for detection and retrieval (i.e., re-identification in this case). For the DC-I-Net, the proposals of each stream are fed into the fully-connected layers for refined detection results. After that, the object features generated by the ROI-Align layers based on the refined detection results are concatenated and fed into another full-connected layer for re-identification. The essential differences between I-Net and DC-I-Net lie in two aspects. 1) For the former, detection and re-identification are treated separately in different layers. 2) For the latter, the two-stage refined objects are used for re-identification rather than the one-stage proposals.

called DC-I-Net, is proposed, as is shown in Fig. 2. The main difference of DC-I-Net from the I-Net lies in that each task is performed on different feature layer, such that the more precise objects after the two-stage object detection are utilized for retrieval. Additionally, in DC-I-Net, we improve the training efficacy of HEP loss by exploiting the class centers of the priority classes, and propose a class center guided HEP, called C^2 HEP loss.

This paper is a substantial extension of our I-Net, where we have made the following new contributions in network architecture and feature loss function:

- Consider the specificity of each task, with the philosophy of divide and conquer, an improved DC-I-Net with new network structure is proposed. The improved structure treats the tasks differently by specially deploying different layered features rather than the shared feature representation. Also, the features for retrieval are from precisely positioned objects instead of coarse proposals, such that the flexibility and plasticity in the co-learning of each task can be improved over the I-Net.
- An improved HEP loss, called C^2 HEP, is proposed for promoting the training efficacy of traditional softmax based cross-entropy loss. The C^2 HEP loss is formulated by using the progressively updated class centers of each class computed with the timely updated input features.
- Extensive experiments on benchmark datasets, such as CUHK-SYSU dataset and PRW dataset for person search and the large-scale WebTattoo dataset for tattoo search, demonstrate that the proposed DC-I-Net achieves another new record based on I-Net and outperforms many other state-of-the-art image search models of task-separately and jointly trained.

The remainder of this paper is organized as follows. Section 2 presents the related work. Section 3 presents the proposed I-Net framework. The new materials of the proposed DC-I-Net is presented in Section 4. The experiments

and results are presented in Section 5. The model analysis and discussions of the proposed models are presented in Section 6. Finally, Section 7 concludes this paper.

2 RELATED WORK

This paper mainly addresses an image search issue together with person search and tattoo search. Consider that a number of *person* related research work has been widely studied in computer vision in recent years, we therefore briefly revisit the closely-related pedestrian detection, person re-identification and some existing person search works.

2.1 Pedestrian Detection

Pedestrian detection is an important branch in object detection. In the early years, traditional object detection methods were designed based on the hand-crafted features and AdaBoost classifiers, such as ACF [18], LDCF [19], Checkerboards [20] and Integral Channels Features (ICF) [21]. In 2005, Dalal and Triggs [22] proposed the Histograms of oriented gradients (HOG) with support vector machine (SVM) classifier which first opened the research of human detection. These traditional methods dominate the field of detection for many years due to their robustness and effectiveness. Motivated by the great success of the convolutional neural networks (CNN), many deep learning based pedestrian detection methods have been developed in recent years. Tian *et al.* [23] jointly optimized pedestrian detection with semantic tasks, including pedestrian attributes and scene attributes. Song *et al.* [5] combined multiple deep networks with one fully-connected layer to improve the detection accuracy. In [24], CNN features extracted by a region proposal network (RPN) [11] are fed into the random forest for pedestrian detection. As a famous two-stage detector, Faster-RCNN [11] generates proposals in the first stage and refines the object for more accurate detection in the second stage, which is trained in an end-to-end manner and achieves state-of-the-art detection performance. In this

paper, Faster-RCNN [11] is taken into account as the object detector in the integrated image search network.

2.2 Person Re-identification

Person re-identification (Re-ID) aims to match a query person (probe) from a set of person candidates (gallery), where the probe and gallery are captured from non-overlapped camera views, which have attracted lots of attention in recent years [25], [26], [27], [28], [29], [30], [31], [32], [33], [34]. Person Re-ID is still an open yet challenging issue for real-world surveillance application, due to the diverse variations of human poses, camera views, backgrounds, illumination, occlusions and resolutions. Earlier person re-identification milestones focus on the feature representation [35], [36], [37], [38], [39] and similarity metric learning [40], [41], [42]. Chen *et al.* [37], decomposed the person re-identification task into the classification and ranking sub-tasks, and jointly optimized them during the training phase. Cheng *et al.* [13] presented a novel multi-channel parts-based model with triplet loss to learn discriminative feature representation. Some of the person re-identification methods in recent years were designed based on mask learning. Xu *et al.* [25] masked out the background of the person and generated a pose-guided score. Song *et al.* [16] introduced a three-stream network with triplet loss. Recently, Sun *et al.* [43] proposed a part-based convolutional baseline (PCB) which fully exploited the local part-level information for feature discrimination. Additionally, in order to remit the stagnation problem of the triplet loss, Chen *et al.* [44] enlarged the three-stream network to quadruplet network such that one more negative pair can be obtained to restrict the positive pairs. In this work, we propose an OLP loss by generating a number of negative samples to restrict the positive pair, which can effectively alleviate the stagnation problem and make end-to-end training much easier.

2.3 Person Search

Person search, as one typical task of image search, have very recently attracted people's attention, which can be divided into two branches: individually trained models and jointly trained models. For the former, the ID-discriminative Embedding (IDE) and Confidence Weighted Similarity (CWS) were firstly proposed by Zheng *et al.* [45] for person search, in which the detector and re-identifier were trained individually. Recently, MSM [46], MGTS [47] and Local Refinement [48], also abandon the end-to-end network structure and deploy two distinct backbones for detection and person re-identification, respectively. After detection, these models convert the person search into a person re-identification (ReID) task. Although these models can work well in person search benchmarks, the intrinsic relationship between detection and re-identification is neglected. Additionally, the individually trained models lose the advantages of low computational cost and user-friendly property of person search. Instinctively, the detector and re-identifier can correspond and adapt to each other during joint training phase.

For the latter, NPSM [49] introduced a LSTM based end-to-end person search method which automatically reduces the region containing the target person from a given image. Yan *et al.* [50] firstly introduced the GCN in person

search for exploring the relation between instances in an image based on the context information and achieved SOTA performance. Xiao *et al.* [51] jointly trained the detection and person re-identification parts during the training phase, in which a classical OIM loss function is introduced for person re-identification. However, the OIM loss only regards the feature learning of Re-ID as a classification problem, which may not well capture the intra-similarity and inter-dissimilarity in feature representation. In contrast, our proposed On-line Pairing (OLP) loss with a Siamese architecture can learn effective similarity metric for discriminative representation [17]. The newly designed multi-task network in this paper allows it to integrate the detection loss, metric loss and classification loss to train simultaneously for more accurate and user-friendly person search.

3 THE PROPOSED INTEGRATED NET (I-NET)

In this section, we will introduce our Integrated Net (I-Net) in details. The I-Net is proposed to jointly handle the detection and person re-identification into an end-to-end framework for user-friendly image search task. The Siamese architecture of I-Net is shown in Fig. 2 (middle). For each iteration, a pair of images containing objects (e.g., persons) of the same identity are fed into the Siamese I-Net. The backbone network is used for preliminary feature representation. After that, two region (pedestrian) proposal networks (RPN) are implemented to get the person proposals in each image, respectively. These proposals are then fed into ROI-pooling layers and the output feature maps are formulated, which is followed by two fully-connected layers for detection task and an extra 256-D \mathcal{L}_2 -normalized feature layer for retrieval (i.e., re-identification) task. Two loss functions, i.e., OLP and HEP losses, are proposed for learning useful features with respect to re-identification.

3.1 Novel Design of Network Structure

Backbone. Our I-Net is specially designed based on the Siamese structure. The backbone of I-Net is based on the VGG16 [52], which has five stacks of convolutional layers, including 2, 2, 3, 3, and 3 convolutional layers for each stack. 4 max-pooling layers are followed for the first four stacks. On the top of *conv5_3* layer, we generate 512 channelled feature maps for predicting the pedestrian proposals. A $512 \times 3 \times 3$ convolutional layer is first added to get the features for computing the pedestrian proposals.

Object Proposals. Similar to the faster RCNN [11], we also associate nine anchors at each feature map. Then, a softmax classifier (cls.) supervised by cross-entropy loss is used to predict whether the anchor is a pedestrian or not, and a Smooth- \mathcal{L}_1 Loss (reg.) is used for bounding box regression. Finally, 128 proposals from each image after the non-maximum suppression (NMS) are obtained, which is generally recognized as the 1st stage detection. Note that, the two branches of the Siamese network are weights shared in our model, and the two RPNs simultaneously generate the object (pedestrian) proposals from each two given input images for subsequent Re-ID task.

Joint Detector and Re-identifier. With the generated proposals from both two RPNs, the ROI pooling layer [53] is

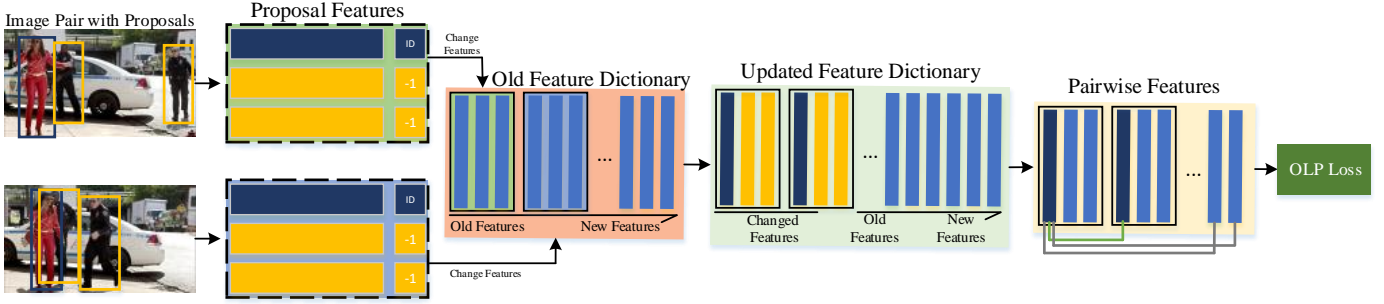


Fig. 3. The flow-chart for computing the OLP loss. The detected proposals include two types of pedestrian, i.e., person with identity (p-w-id) in dark blue box and person without identity (p-w/o-id) in yellow box. The yellow box is labeled as -1. These proposal features are stored in the feature dictionary, which is used to construct the pair-wise features including positive pairs (green lines) and negative pairs (gray lines).

integrated into our I-Net to pool the features from the convolution layer. The pooled features from both two branches are then fed into the two fully-connected layers with 4096 neurons. In order to remove the false positives from the pedestrian proposals, a two-class (binary) softmax classifier (cls.) for differentiating person from others (non-persons) is trained by cross-entropy loss. Note that, for general image search task, a multi-class softmax classifier will be trained for general object detection. The Smooth- \mathcal{L}_1 loss is used to refine the locations of bounding boxes. This is generally recognized as the 2nd stage detection similar to Faster-RCNN. Then a pair of 256-D \mathcal{L}_2 -normalized features for each pair-wise images generated by an extra fully-connected layer are fed into the on-line pairing (OLP) loss and hard example priority (HEP) loss for training the re-identification module, and the details of OLP and HEP will be presented in the following sections. With these two-stage detection losses (cls. vs. reg.) and re-identification losses (OLP vs. HEP), the proposed I-Net can be jointly trained for simultaneous person detection and re-identification in an end-to-end Siamese network architecture, as is shown in Fig. 2.

Overall Merits. There are two clear merits of the Siamese structure. On one hand, image search is actually an image matching problem between a query image (probe) and each gallery image, which means that the two-stream I-Net can appropriately match two given images, by training the metric loss and classification loss simultaneously. On the other hand, with the pair-wise images as inputs, the number of samples is increased, which is beneficial to the training efficacy and robustness of the detector and re-identifier.

3.2 On-line Pairing Loss (OLP)

The OLP loss is proposed for Re-ID task by learning discriminative metric, with the consideration of two aspects:

- First, the features for person re-identification should be restricted by the loss function during the training phase, such that the features of the same identity should have smaller distance (intra-similarity) and the features of different identity should have larger distance (inter-dissimilarity).
- Second, due to the insufficient input image numbers and lack of objects in each image, the stagnation problem of the traditional metric loss (e.g., triplet loss) easily happens because of many easy pairs but

few identities (few-shot), which seriously prevents the model from being effectively trained.

With the above considerations, in the model we instinctively deploy a dynamic feature dictionary to store the proposal features, such that more negative pairs can be generated together with the positive pairs drawn from more identities. As a result, an effective metric with these positive and negative pairs in a Siamese structure can be learned. The reason is that the condition of the loss function is much harder to be satisfied because of many more pairs and identities shot, and the stagnation problem in training can then be remitted.

Specifically, for each iteration, the 1st stage detector provides 128 bounding boxes (proposals) in each image. In the person search datasets, without considering the proposals of backgrounds, there are two types of pedestrian bounding boxes, i.e., persons with identity information (p-w-id) and persons without identity information (p-w/o-id). The detailed flow chart of the proposed OLP loss is described in Fig. 3, in which the p-w-id and p-w/o-id are represented by dark blue and yellow bounding boxes, respectively. We observe that an online feature dictionary of fixed size is deployed to store the generated features together with their identity label (i.e., -1 for the p-w/o-id). In this work, the number of features stored in the dictionary (i.e., dictionary size) is set as 40 times the number of bounding boxes generated by the detector. When the dictionary is filled with features during training phase, the oldest features in the dictionary will be replaced with new ones.

Formulation. In order to minimize the discrepancy between the features of the same identity and simultaneously maximize the discrepancy between the features of different identity, we have the following notations for formulation of the OLP loss. Suppose the group of proposals for loss computation to be $(\mathbf{p}_1, \mathbf{p}_2, \mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_K)$, where $(\mathbf{p}_1, \mathbf{p}_2)$ stand for bounding boxes from the same identity generated by the model in forward propagation, $(\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_K)$ are the features stored in the dictionary labeled as negative samples, and K means the number of negative samples for the i^{th} subgroup. For each pairing group, we tend to formulate two symmetrical subgroups by taking \mathbf{p}_1 and \mathbf{p}_2 as anchor, alternatively. For example, when \mathbf{p}_1 is regarded as anchor, then $(\mathbf{p}_1, \mathbf{p}_2)$ denotes the positive pair, while $(\mathbf{p}_1, \mathbf{n}_1), (\mathbf{p}_1, \mathbf{n}_2), \dots, (\mathbf{p}_1, \mathbf{n}_K)$ represent negative pairs. Alternatively, when \mathbf{p}_2 is regarded as anchor, then $(\mathbf{p}_2, \mathbf{p}_1)$ denotes the positive pair,

while $(\mathbf{p}_2, \mathbf{n}_1), (\mathbf{p}_2, \mathbf{n}_2), \dots, (\mathbf{p}_2, \mathbf{n}_K)$ represent negative pairs. Obviously, the OLP loss function generates more negative pairs to restrict the positive pair, which is able to remit the stagnation problem.

Suppose that we get m subgroups in one iteration, and $\mathbf{x}_a^i, \mathbf{x}_p^i, (\mathbf{x}_{n_1}^i, \mathbf{x}_{n_2}^i, \dots, \mathbf{x}_{n_K}^i)$ stand for the anchor, positive and negative features of i^{th} subgroup, respectively. Then, the proposed OLP loss function is represented as follows.

$$\mathcal{L}_{OLP} = -\frac{1}{m} \sum_{i=1}^m \log \frac{e^{d(\mathbf{x}_a^i, \mathbf{x}_p^i)}}{e^{d(\mathbf{x}_a^i, \mathbf{x}_p^i)} + \sum_{k=1}^K e^{d(\mathbf{x}_a^i, \mathbf{x}_{n_k}^i)}} \quad (1)$$

where $d(\cdot)$ stands for the cosine similarity between two features. It is worth noting that since our features are \mathcal{L}_2 -normalized, the cosine similarity can be easily computed by the inner product of each two feature vectors. From Eq. (1), we observe that a softmax guided cross-entropy loss works as a metric loss, and the summation of all distances is set as the denominator so that the distance of the positive pair can compare to all negative pairs in each subgroup.

In gradient computation, we only calculate the deviation with respect to the anchor feature. Then, the deviation of the OLP loss function with respect to \mathbf{x}_a^i for the i^{th} subgroup can be calculated as:

$$\frac{\partial \mathcal{L}_{OLP}}{\partial \mathbf{x}_a^i} = (q^i - 1)\mathbf{x}_p^i + \sum_{k=1}^K (\hat{q}_k^i \mathbf{x}_{n_k}^i) \quad (2)$$

where q^i and \hat{q}_k^i are expressed as follows.

$$q^i = \frac{e^{d(\mathbf{x}_a^i, \mathbf{x}_p^i)}}{e^{d(\mathbf{x}_a^i, \mathbf{x}_p^i)} + \sum_{k=1}^K e^{d(\mathbf{x}_a^i, \mathbf{x}_{n_k}^i)}} \quad (3)$$

$$\hat{q}_k^i = \frac{e^{d(\mathbf{x}_a^i, \mathbf{x}_{n_k}^i)}}{e^{d(\mathbf{x}_a^i, \mathbf{x}_p^i)} + \sum_{k=1}^K e^{d(\mathbf{x}_a^i, \mathbf{x}_{n_k}^i)}}, k = 1, \dots, K \quad (4)$$

In summary, as is shown in Fig. 4 and Eq. (1), the proposed OLP loss can be implemented as follows:

- 1) The features of each two input images are collected. The features $(\mathbf{p}_1, \mathbf{p}_2)$ from the images of the same identity are constructed as positive pairs.
- 2) For each positive pair $(\mathbf{p}_1, \mathbf{p}_2)$, \mathbf{p}_1 and \mathbf{p}_2 are set as the anchor, alternatively. The features $(\mathbf{n}_1, \mathbf{n}_2, \dots, \mathbf{n}_K)$ stored in the feature dictionary are paired with the anchor to construct negative pairs.
- 3) Compute the OLP loss by using Eq. (1) and its gradient by using Eq. (2) for gradient back-propagation optimization.
- 4) Store the input features to progressively update the feature dictionary.

Obviously, by optimizing the proposed OLP metric loss function, the cosine similarity $d(\mathbf{x}_a^i, \mathbf{x}_p^i)$ between the features of the same identity (intra-similarity) is maximized, while the cosine similarity $d(\mathbf{x}_a^i, \mathbf{x}_n^i)$ of different identities (inter-similarity) is minimized. Moreover, with thousands of features progressively stored in the feature dictionary, a number of negative pairs can be generated which also effectively remit the stagnation of I-Net model training. Note that our OLP is scenario driven and tailored for the multi-task integration network, which is essentially different from the pure metric learning loss, such as the N-pair loss

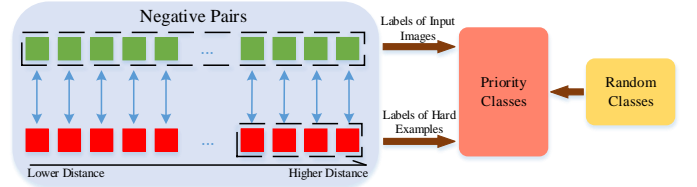


Fig. 4. The protocol for selecting the priority classes of hard examples for computing the HEP loss. First, the person proposals (bounding boxes) with identities (i.e., ground-truth labels) are marked. Second, the negative pairs with the largest cosine distances were selected as hard examples, which are denoted as priority classes. Finally, if the pool of priority classes is not yet filled, some random classes are selected to fill the pool, which are used to compute the HEP loss.

[54] that relies on enough pair-wise training samples from enough classes. The existing metric losses cannot easily be deployed in our network for feature learning due to that the number of input training samples is too small.

3.3 Hard Example Priority Loss (HEP)

The OLP loss function enables the cosine distance of positive pairs to be smaller and that of negative pairs to be larger, which does not directly regress the identity labels in the loss function. Additionally, the traditional softmax based cross-entropy loss for classifier training does not consider the degree of difficulty of the examples in the data. With the above considerations, we further propose a hard example priority (HEP) loss function, which aims to regress the identity labels with high priority. A high priority of some identity label means that the identity is hard to be classified, and will be selected for label regression and loss computation.

Suppose that there are C identities. The HEP loss function aims to classify the person proposals with identity (i.e., p-w-id) into $C + 1$ classes containing an extra background class. In order to calculate the HEP loss, the person proposals with high Interaction-over-Union (IOU) between the bounding boxes and the ground truth are used. Then, by computing the cosine similarity between positive pairs and negative pairs via the OLP loss as described in Section 3.2, we can determine the top r maximum distances of the negative pairs. The examples with maximum distances of negative pairs (i.e., high inter-similarity) are recognized as hard examples of priority classes. In order to keep the total number of priority classes fixed, we also randomly select an uncertain number of classes from the remaining categories of non-priority classes. As a result, totally $T(T \ll C + 1)$ classes are selected to compute the HEP loss. In summary, suppose that the selected hard categories are stored in the pool \mathcal{P} , as is shown in Fig. 4, the protocol for selecting the hard categories is presented as follows:

- 1) The label indexes of each input image pairs with identities are first determined to ensure the ground-truth classes.
- 2) For each subgroup (described in Section 3.2), the label indexes of negative samples from the top r negative pairs with the maximum distances are stored in the priority classes pool \mathcal{P} , which enables the priority classes of hard examples to be focused.

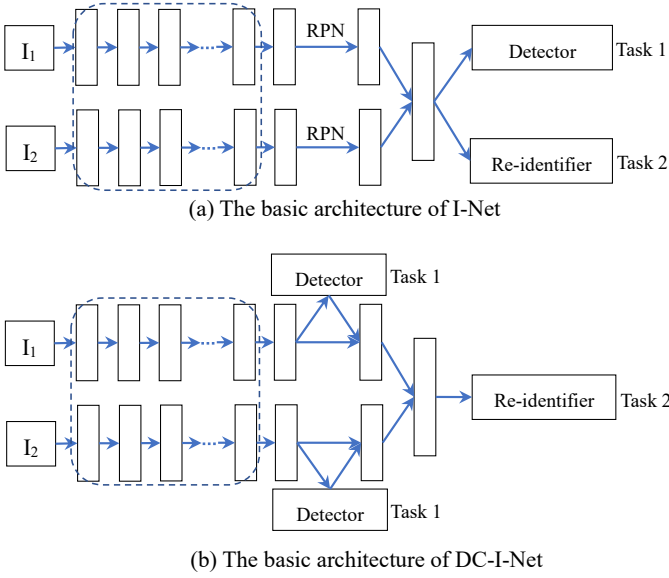


Fig. 5. Motivation of DC-I-Net and its differences from I-Net that 1) the detector is deployed in front of the re-identifier in different layers and 2) the refined objects instead of the coarse proposals from RPN are used to train the re-identifier.

- 3) If the size of pool \mathcal{P} is still smaller than the preset T , we randomly select several classes to fill the pool.

Finally, with the traditional softmax based cross-entropy loss and the selected priority classes, the proposed HEP loss function is formulated as:

$$\mathcal{L}_{HEP} = -\frac{1}{n} \sum_{i=1}^n \sum_{j \in \mathcal{P}} \mathbf{1}(\text{label} = j) \log \frac{e^{s_j^i}}{\sum_{t=1}^T e^{s_t^i}} \quad (5)$$

where s_j^i stands for the i -th proposal's score from the classifier and j stands for the j -th class. In the loss function, only the selected categories are used for the loss computation, such that, the loss function focus on the hard category.

3.4 Overall Loss of I-Net

I-Net is an end-to-end model which has integrated the detection and re-identification jointly for training. Therefore, the losses are composed of two parts: detection loss (\mathcal{L}_{Det}) and re-identification loss (\mathcal{L}_{OLP} and \mathcal{L}_{HEP}), which are represented as follows.

$$\mathcal{L}_{I-Net} = \mathcal{L}_{Det} + \alpha \mathcal{L}_{OLP} + \beta \mathcal{L}_{HEP} \quad (6)$$

where α and β represent the trade-off parameters.

Note that the detection loss \mathcal{L}_{Det} is following the traditional Faster-RCNN detector, which includes two stages and each stage includes a softmax based cross-entropy classification loss (cls.) and smooth- \mathcal{L}_1 based regression loss (reg.) for bounding box prediction. The training of model (6) is end-to-end for user-friendly image search task.

4 DIVIDE AND CONQUER INTEGRATED NETWORK (DC-I-NET)

Motivation. With the philosophy of divide and conquer, this section presents an improved I-Net, i.e., DC-I-Net, for

dealing with several flaws of I-Net in network architecture and training loss. The difference in architecture between the I-Net and DC-I-Net is clearly shown in Fig. 5, which, specifically, is motivated and improved from the following three important aspects:

- *Task specialization in architecture for detection and Re-ID.* We propose a novel network structure over the I-Net, in which the features for detection and re-identification are extracted from different layers rather than the shared layer as I-Net does, by deploying the detector in front of the re-identifier. The reason is that the features for the detection task should focus on the discrimination between the foregrounds (objects) and backgrounds, while the features for re-identification should focus on the discrimination among the foregrounds (e.g., persons of different identities). The novel network architecture deployed with detector and re-identifier is trained end-to-end, such that the image search task becomes more user-friendly and practical.
- *Refined object proposals for the metric loss.* In I-Net, the features from the coarse object proposals (the 1st stage detection) are used for re-identification loss optimization, which may degrade the final retrieval performance due to the inaccurate proposals. Therefore, in DC-I-Net, the refined objects with ROI-Align (the 2nd stage detection) are used for better training the re-identification metric loss.
- *Easy training of the classification loss.* As mentioned in Section 3.2, due to the small input number of images during the training phase, a very few number of identities contribute to the model training such that the model has to experience thousands of epoches for probably seeing all labeled identities. This naturally leads to the hard training of the HEP loss. In order to address this issue, we further propose a class center guided hard example priority (C^2 HEP) loss by fully exploiting the updated input features to compute the class centers. As a result, the identity discrimination of features is much improved.

From the perspectives of network architecture and loss function, these new improvements of DC-I-Net over I-Net are presented in the following subsections in detail.

4.1 Improved Network Architecture Over I-Net

The backbone of DC-I-Net is the same as I-Net, as shown in Fig. 5 (the dashed box). The detailed network structure of DC-I-Net is shown in Fig. 2 (right), which is different from the I-Net shown in Fig. 2 (middle) that 1) the task specialization for detection and Re-ID is well considered by using features from different layers, 2) the refined objects from the 2nd stage detector are generated by using ROI-Align module for training the metric loss, and 3) a class-center guided hard example priority (C^2 HEP) loss is proposed for easy training of the identity classification loss.

Detector. In the DC-I-Net, the features for detection and person re-identification tasks are extracted from different layers. After the two-stage detection supervised by classification loss (cls.) and regression loss (reg.), the detection of accurate bounding boxes (i.e. objects) is completed.

Re-identifier. After the two-stage detection, the coordinates of the refined bounding boxes are fed into the ROI-Align layers to compute the features of refined object proposals for person re-identification (person search) or object retrieval (image search). The pooled feature maps have a size of 7×14 for person search task, which has a similar aspect ratio to the bounding boxes of a person. The feature maps are then fed into the fully-connected layers to learn the feature vector representation for person re-identification. Finally, 256-D \mathcal{L}_2 -normalized features are generated for the object proposals by an extra fully-connected layer, which are then fed into the OLP loss and $\mathcal{C}^2\text{HEP}$ loss for formal training of the re-identification module.

4.2 Class Center Guided HEP Loss ($\mathcal{C}^2\text{HEP}$)

The class-center guided HEP loss ($\mathcal{C}^2\text{HEP}$) is proposed to improve the hard training problem of HEP in I-Net, which is an inherent problem of person search task caused by the insufficient number of identities of each iteration. Therefore, we exploit the features extracted from the input images to compute the class center of each category. Consider that the cosine similarity between each sample and the class centers can clearly reflect the probability of the sample belonging to a class, we propose to feed the cosine similarity into the softmax function of HEP and formulate the $\mathcal{C}^2\text{HEP}$. For convenience, we design a class center dictionary indicated by their ground-truth label. In every iteration, the stored class center for class j is updated by new features, formulated as

$$\mathbf{c}_{new}^j = \phi \cdot \mathbf{c}_{old}^j + (1 - \phi) \cdot \mathbf{x}^{(j)} \quad (7)$$

where \mathbf{c}_{old}^j is the old center of class j , $\mathbf{x}^{(j)}$ is the input feature of class j , and ϕ ($0 < \phi < 1$) is a hyper-parameter, which is set as 0.5 in our implementation. It is worth noting that the class center dictionary is different from the feature dictionary in the OLP loss. In order to fully explore the information of identities/categories, each labeled object (e.g., person) in the dataset has been assigned a class center stored in the class center dictionary.

Suppose that the feature of the i^{th} object fed into the $\mathcal{C}^2\text{HEP}$ loss function is \mathbf{x}_i , the class center dictionary is defined as \mathcal{S} , where $\mathcal{S} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_C\}$ and C is the number of the identities (categories) annotated in the dataset. Generally, the feature more probably belongs to the class that has the smallest cosine distance between the class center and the feature. Based on the softmax function, we define the probability of \mathbf{x}_i belonging to the class j is defined as:

$$p_j = \frac{e^{\lambda d(\mathbf{x}_i, \mathbf{c}_j)}}{\sum_{c \in \mathcal{P}} e^{\lambda d(\mathbf{x}_i, \mathbf{c}_c)}} \quad (8)$$

where the λ is a hyper-parameter, which is set as 10 in implementation and \mathcal{P} is the pool of selected priority classes, which is presented in Section 3.3. From Eq. (8), we know that the highest probability p_j can be achieved when a given sample \mathbf{x}_i belonging to the class j has the smallest cosine distance to the class center \mathbf{c}_j .

Suppose that the probability of a sample belonging to their corresponding ground truth label is represented as $\{p^1, p^2, \dots, p^n\}$, where n is the number of samples in one iteration. Under the assumption of independently and identically distribution of the learned features, it is rational to

maximize the likelihood function: $L = \prod_i^n p^i$ for model training. In order to train the deep learning model, we minimise the negative log-likelihood function: $-\log L = -\log \prod_i^n (p^i)$. With the negative log-likelihood function, the proposed $\mathcal{C}^2\text{HEP}$ loss function can be written as:

$$\mathcal{L}_{\mathcal{C}^2\text{HEP}} = -\frac{1}{n} \sum_{i=1}^n \sum_{j \in \mathcal{P}} \mathbf{1}(\text{label} = j) \log \frac{e^{\lambda d(\mathbf{x}_i, \mathbf{c}_j)}}{\sum_{l \in \mathcal{P}} e^{\lambda d(\mathbf{x}_i, \mathbf{c}_l)}} \quad (9)$$

From Eq. (9), we know that minimizing the loss function can effectively constrain the cosine similarity $d(\cdot)$ between each feature and its corresponding class center to be larger. Therefore, the identity discrimination is guaranteed. Note that the proposed scenario-driven $\mathcal{C}^2\text{HEP}$ loss is tailored for the detection and re-identification integration network, which is explicitly and implicitly different from the existing class-center based feature learning loss, such as center loss [55] and prototype loss [56], that only relies on the available training samples. However, in person search scenario, the number of training samples in each iteration is very small (i.e. 2 samples for each iteration in our model), and the center loss and prototype loss will encounter stagnation problem. Therefore, a class-center dictionary is deployed in $\mathcal{C}^2\text{HEP}$ for effectively alleviating the training stagnation problem and improving the training efficiency. The $\mathcal{C}^2\text{HEP}$ loss can be recognized as a seamless connection to the OLP loss for feature discrimination.

4.3 Overall Loss of DC-I-Net

The DC-I-Net is also an end-to-end model for user-friendly image search task. Similar to I-Net, the losses of DC-I-Net consist of the detection loss and re-identification loss. Specifically, the detection loss \mathcal{L}_{Det} is following the traditional Faster-RCNN with two-stage detection, and each stage refers to the softmax based cross-entropy classification loss (cls.) and smooth- \mathcal{L}_1 regression loss (reg.) for bounding box prediction. Since the two streams of the Siamese network share the same parameters for the detection, the two input images are used to train the detector simultaneously.

Besides, the re-identification loss consists of the proposed OLP metric loss and $\mathcal{C}^2\text{HEP}$ identity classification losses for discriminative feature representation. So the overall loss of the DC-I-Net is presented as follows:

$$\mathcal{L}_{DC-I-Net} = \mathcal{L}_{Det} + \alpha \mathcal{L}_{OLP} + \beta \mathcal{L}_{\mathcal{C}^2\text{HEP}} \quad (10)$$

where α and β stands for the trade-off parameters, and \mathcal{L}_{OLP} is given in Eq. (1). By using the mini-batch SGD optimization, our model can be trained end-to-end for person search. In summary, four steps are involved during the training in every iteration:

- 1) *Detection loss computation.* The pair-wise input images are fed into the Siamese structure for detection first, and the loss of detection for each image is computed.
- 2) *OLP loss computation.* The detected objects (feature maps) are fed into the ROI-Align layer to get the features for re-identification task. The features are paired, and the distances of positive and negative pairs are calculated. Then, the OLP loss is computed with dynamic update of the feature dictionary.

- 3) *C²HEP loss computation.* The distances of positive and negative pairs are fed into the C²HEP loss with selected priority classes. After computation of the loss, the class centers are progressively updated via the input new features.
- 4) *Gradient computation.* Based on the computation of all losses, the model is optimized with SGD, until convergence.

5 EXPERIMENTS

To evaluate the effectiveness of our approaches, we conduct massive experiments on three benchmark datasets, including CUHK-SYSU dataset [57], PRW dataset [45] and Webtattoo dataset [58], for image search tasks. The first two datasets focus on the person image search (i.e., person search), which refers to joint object (person) detection and person re-identification tasks in our models. The third dataset serves for tattoo image search (i.e., tattoo search), which refers to joint object (tattoo) detection and image retrieval tasks in our models. In this section, the experimental setup and experimental results for each dataset are presented.

5.1 Experimental Setup

5.1.1 Implementation Details

The proposed I-Net and DC-I-Net are implemented on Caffe [59] and py-faster-rcnn [11] platform for model training and evaluation. The VGG-16 [52] is used as the backbone network of our models and the pre-trained model in [57] is taken into account for network parameters initialization. The first two stacks of convolutional layers are frozen during the training of our models. The two branches of the Siamese network share the same parameters for both initialization and training. The RPN part of each branch generates 128 proposals for each image, and the proposals labeled as background are not useful and therefore dropped for object retrieval task. In both I-Net and DC-I-Net, the trade-off parameters α and β are set as 1. The learning rate is initialized to 0.001, and drops to 0.0001 after 40k iterations. Totally, 70k iterations are set to enable convergence.

5.1.2 CUHK-SYSU Dataset

The CUHK-SYSU dataset [57] is a large dataset for person search, which contains 18184 images from the hand-held cameras and movie snapshots with large variations in viewpoint, lighting, resolution, *etc.* From the annotations, there are 8432 different person identities and 96143 bounding boxes. Each labeled person has at least two images from different viewpoints. For the training/testing split, the developer of this dataset provided 11206 images of 5532 identities for training and 6978 images of 2900 identities for the test. Specifically, we follow the same experimental protocols as [57] for fair comparison.

5.1.3 PRW Dataset

The PRW dataset [45] is drawn from a 10 hours of video captured by six cameras, in which five of them are 1080 × 1920 HD and the remaining one is 576 × 720 SD. Totally 11816

frames are manually annotated and results in 43110 pedestrian bounding boxes, in which 34304 pedestrians are annotated by 932 IDs. For the training/testing split, the PRW dataset provides 5134 frames of 482 labeled identities for training and 6112 frames of 450 labeled identities for testing. The task for this dataset allows the model to search a query target person (probe) from the whole testing set (gallery), which remains to be a challenging problem.

5.1.4 Webtattoo Dataset

The Webtattoo dataset [58] was presented in different viewpoints and illuminations, which consists of three parts: (i) the first part is a combination of three small-scale (less than 10K) tattoo datasets, such as Tatt-C [60], Flickr [61] and DeMSI [62]. (ii) The second part is a collection over 300K distracter tattoo images from the Internet. (iii) The third part is the 300 tattoo sketches drawn by volunteers. In this Webtattoo dataset, three tasks including the detection, tattoo search and sketch based tattoo search are deployed. In this paper, we focus on the joint tattoo detection and image search. Specifically, 1428 images of 400 tattoo classes are used for model training. For comparing the detection performance of different models, 755 images from 200 tattoo classes with ground-truth bounding boxes are used. For comparing the search (retrieval) performance of different models, the query set containing 200 images (one image per tattoo class) is used to search the images from a gallery set containing 355 tattoo images.

5.2 Experiments on CUHK-SYSU Dataset

5.2.1 Compared Methods

Baselines: Separated Detection and Re-ID Models. In this section, we perform the experiments on the CUHK-SYSU dataset to investigate the effectiveness of our models. Consider that the proposed models in this dataset aim to jointly learning pedestrian detection and person re-identification, we therefore select three pedestrian detection methods and five person re-id approaches for baseline comparisons, which then result in 15 baselines for person search task. Specifically, three baseline detection methods, CCF [63], Faster-RCNN [11] with Resnet50 [64] and ACF [18], are used for detecting pedestrians. Besides, we also use the ground truth bounding boxes of the test set as the upper bound of the detector's performance. For the baseline re-identification methods, we evaluate several famous re-id feature representation methods including DenseSIFT-ColorHist (DSIFT) [65], Bag of Words (BoW) [66], Local Maximal Occurrence (LOMO) [40] and ID-Net(The re-identification part of OIM [51]). The metric learning methods, i.e. KISSME [41] and XQDA [40] together with these feature representation are used for Re-ID. These separated detection and Re-ID methods are combined for person search, which are therefore treated as baselines in comparisons.

State-of-the Art (SOTA): Joint Detection and Re-ID Models. To the best of our knowledge, there is only a few work on the joint training of detector and re-identifier for person search task, such as the OIM model [51], the end-to-end model (initialized model) [57], NPSM [49], IAN [67], RCAA [68] and Context Graph [50]. Therefore, these end-to-end person

TABLE 1
Comparisons of baselines, SOTA methods and our models on the CUHK-SYSU dataset

Detector	Re-id Method	mAP(%)	Top-1(%)
ACF	DSIFT [65]+Euclidean	21.7	25.9
	DISFT [65]+KISSME [41]	32.3	38.1
	BOW [66]+KISSME [41]	42.4	48.4
	LOMO [40]+XQDA [40]	55.5	63.1
	IDNet [51]	56.5	63.0
CCF	DSIFT [65]+Euclidean	11.3	11.7
	DISFT [65]+KISSME [41]	13.4	13.9
	BOW [66]+KISSME [41]	26.9	29.3
	LOMO [40]+XQDA [40]	41.2	46.4
	IDNet [51]	50.9	57.1
CNN	DSIFT [65]+Euclidean	34.5	39.4
	DISFT [65]+KISSME [41]	47.8	53.6
	BOW [66]+KISSME [41]	56.9	62.3
	LOMO [40]+XQDA [40]	68.9	74.1
	IDNet [51]	68.6	74.8
GT	DSIFT [65]+Euclidean	41.1	45.9
	DISFT [65]+KISSME [41]	56.2	61.9
	BOW [66]+KISSME [41]	62.5	67.2
	LOMO [40]+XQDA [40]	72.4	76.7
	IDNet [51]	73.1	78.3
End-to-End(Initialized model) [57]	OIM [51]	55.7	62.7
	IAN [67]	75.5	78.7
	NPSM [49]	76.3	80.1
	RCAA [68]	77.9	81.2
	CNN _v +MGTS [47]	79.3	81.3
	I-Net	83.0	83.7
	Context Graph [50]	79.5	81.5
	DC-I-Net(VGG16)	84.1	86.5
	DC-I-Net(Resnet50)	83.7	85.8
		86.2	86.5

search methods are selected as the SOTA competitor of our I-Net and DC-I-Net. Additionally, the CNN_v+MGTS [47] which trains the detector and person re-identifier separately is also compared, because of their excellent performance. All the experiments are following the same experimental protocols for fair comparisons and the gallery size is set as 100. Note that the I-Net is implemented with VGG16 while the DC-I-Net is implemented with both VGG16 and Resnet50 in the experiment, because almost all the compared deep models are based on Resnet50 backbone.

5.2.2 Experimental Results

In the experiments, the top-1 accuracy and the mAP (mean average precision) are computed for evaluating the performance of person search. Specifically, the results of person search are shown in Table 1, from which we can see that the proposed DC-I-Net with Resnet50 achieves a top-1 accuracy of 86.5% and mAP of 86.2%, and outperforms all the compared methods including the SOTA context graph [50] (2.1% in mAP) that firstly deploys GCN inside. It is worthy noting that, with VGG16 backbone the proposed DC-I-Net outperforms the SOTA end-to-end person search model (i.e., I-Net) by 4.3% and 4.2% in top-1 accuracy and mAP, respectively. From Table 1, we can also observe that the SOTA end-to-end person search methods, such as OIM, NPSM, OP-I-Net, RCAA, and IAN always outperform the traditional person search methods that train the detector and person re-identifier separately rather than a joint way. It is also noteworthy that even with the ground-truth bounding boxes of pedestrians, the traditional re-identification

methods have still shown inferior results compared to the end-to-end detection and Re-ID methods. This demonstrates that the joint training of both detection and re-identification modules is effective and necessary.

Additionally, the recently proposed CNN_v+MGTS [47] is a specially designed model for person search task, which actually trains the detector and re-identifier individually. From the results, we observe that CNN_v+MGTS outperforms other compared separately or jointly trained models. However, it is still inferior to our DC-I-Net by 3.2% and 2.8% in top-1 accuracy and mAP, which further proves our perspective that the joint training of multiple tasks is beneficial to the between-task collaboration and improving the final performance.

5.2.3 Evaluation Remarks

Our model learns the feature representation of person re-identification based on both verification (OLP) and classification (C²HEP) loss functions. Benefited from the Siamese structure and the feature dictionary, our model can easily generate a number of positive pairs and negative pairs to train the OLP metric loss for Re-ID task. Also, the proposed C²HEP loss function can well exploit the progressively updated class centers of each category such that the input features can be trained with better identity classification capability. With these advantages, the proposed DC-I-Net model outperforms the earliest OIM [51] in 2017 for person search by 8.2% and 7.1% in mAP and top-1 accuracy, respectively. Additionally, our DC-I-Net also outperforms the very recent recent RCAA [68] by 4.4% and 4.5% in mAP and top-1 accuracy, respectively. The proposed method also outperforms the single-task work, i.e. CNN_v+MGTS [47], in which each task is separately trained for person search. This demonstrates that joint training of multiple tasks is probable to outperform that of separated training of each task.

5.2.4 Visualization Results of Person Search

Consider that OIM [51] published in CVPR 2017 is the first work for person search and it also deployed a dictionary in the model, we therefore present the visualization results of the DC-I-Net and the OIM in Fig. 6 for comparison, in which eight queries are given and the person search results of top-1, top-2 and top-3 are shown. From the images, the difficulty of person search is observed because of the illuminations, resolutions and dense crowd. In the last two rows, the OIM fails to discover the target person in the top-3 results, while our DC-I-Net can still locate and match all the target persons correctly. From the visualization results, our model manifests better performance and robustness.

However, in some extreme conditions, our model also encounters false alarms (incorrect matches) as is shown in Fig. 7. The extreme conditions are specified as follows. In Fig. 7 (a) and (d), the target persons stay in the crowds and overlapped with each other, which remains a difficult problem in person detection and Re-ID. In Fig. 7 (d) and (e), the person search is affected by illumination, an inherent factor, and the pose variations (b). The similar clothes in Fig. 7 (c), (e) and (f) can easily cause false alarms because of the inter-similarity. Since persons are not always clearly presented in real-world applications, the person search task still faces a challenge.

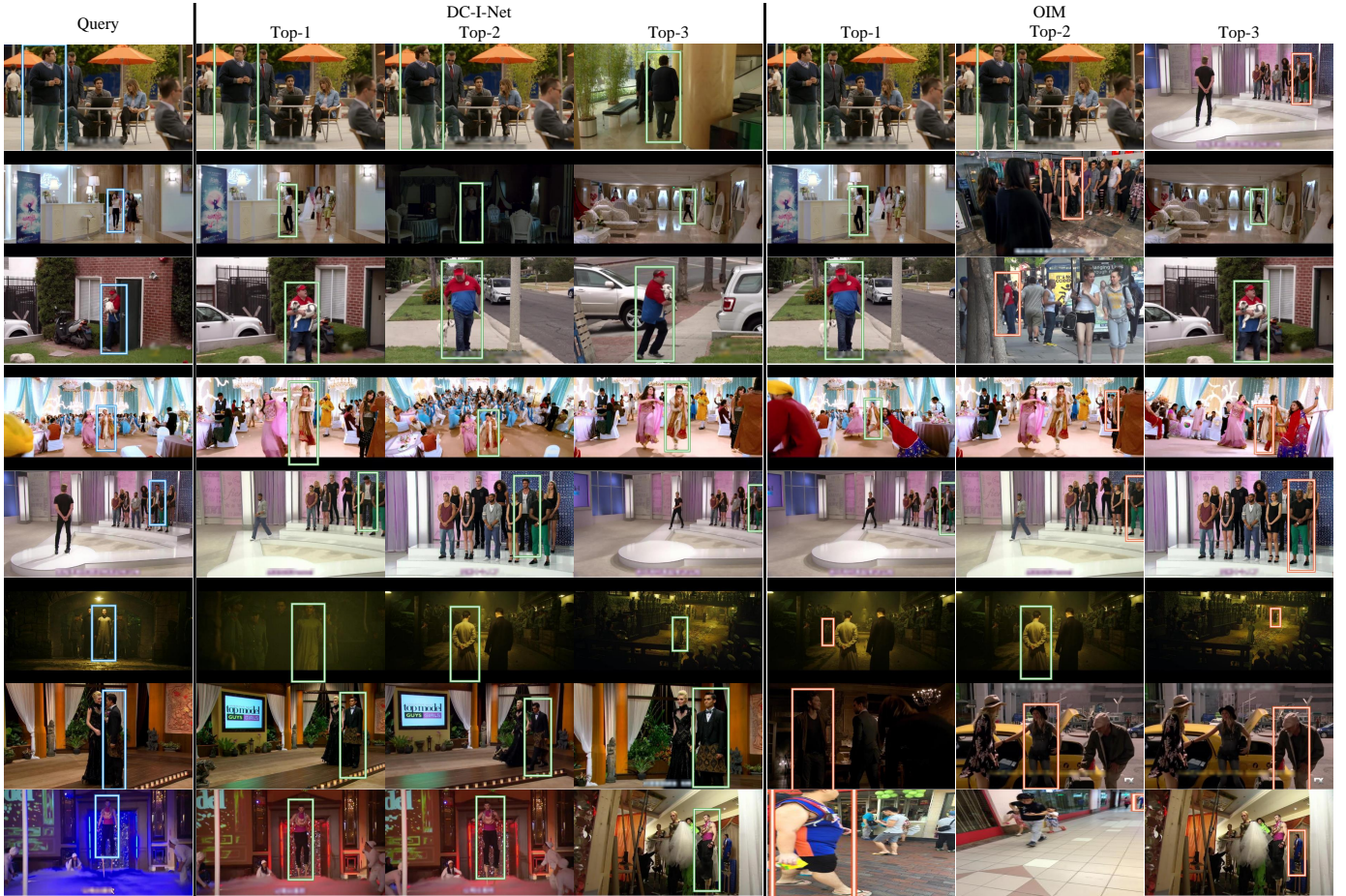


Fig. 6. Retrieval of our DC-I-Net model (middle) and OIM (right) with the CUHK-SYSU dataset by given eight queries (left). The top 3 images with respect to the highest similarity scores are shown. The blue boxes represent the target query person (probe), the green boxes mean correct matches, and the red boxes mean incorrect matches.

We further compare our DC-I-Net with the OIM and our previous I-Net on the CUHK-SYSU dataset from the performance of person search and pedestrian detection by presenting the Precision-Recall (P-R) curves in Fig. 8 (left), respectively. The first row in Fig. 8 denotes the P-R curves of person search performance and the second row represents the P-R curves of pedestrian detection performance. From the P-R curves, we see that the proposed DC-I-Net is much superior to other two closely-related state-of-the-art models for person search. Additionally, as shown in the second row of Fig. 8 (left), the detection of the proposed DC-I-Net is comparable to the individually trained Faster-RCNN [11], which represents the SOTA detection performance. Therefore, the joint multi-task integration model is expected to have better performance beyond the single task of detection by handling multiple tasks appropriately.

5.3 Experiments on PRW Dataset

5.3.1 Compared Methods

Similar to the CUHK-SYSU datasets, for the benchmark PRW dataset [45], we compare our DC-I-Net and I-Net with the SOTA methods for end-to-end person search such as OIM [51] and NPSM [49] and the baseline methods for person search with separate detection and re-id methods. Note that the DC-I-Net is trained with both VGG16

TABLE 2
Comparisons of baselines, SOTA methods and our models on the PRW dataset

Methods	mAP(%)	Top-1(%)
DPM [69]+BOW [66]	9.7	31.1
DPM [69]+IDE _{det} [45]	18.8	45.9
DPM-Alex+LOMO+XQDA [40]	13.0	34.1
DPM-Alex+IDE _{det} [45]	20.3	47.4
DPM-Alex+IDE _{det} + CWS [45]	20.5	48.3
ACF [18]+LOMO+XQDA [40]	10.5	30.9
ACF [18]+IDE _{det} [45]	17.5	43.8
ACF-Alex+LOMO+XQDA [40]	10.3	30.6
ACF-Alex+IDE _{det} [45]	17.5	43.6
ACF-Alex+IDE _{det} + CWS [45]	17.8	45.2
LDCF [19]+BOW [66]	9.1	29.8
LDCF [19]+LOMO+XQDA [40]	11.0	31.1
LDCF [19]+IDE _{det} [45]	18.3	44.6
LDCF [19]+IDE _{det} +CWS [45]	18.3	45.5
OIM [51]	21.3	49.9
NPSM [49]	24.2	53.1
I-Net	25.6	48.7
DC-I-Net(VGG16)	30.4	53.3
DC-I-Net(Resnet50)	31.8	55.1

and Res50. Specifically, for separated detection, the DPM based [69], ACF based [18] and LDCF based [19] methods and their RCNN versions are considered. For separated

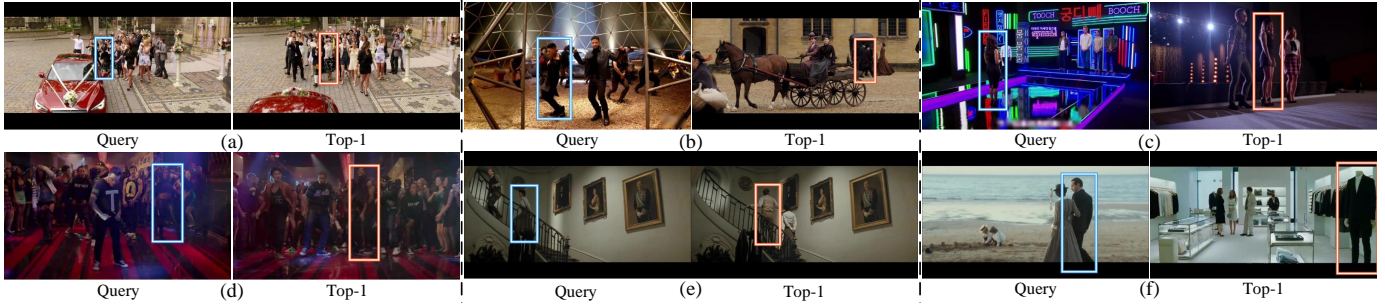


Fig. 7. Some failures on Top-1 of our DC-I-Net model. The impacts from many real-world factors are shown, including the crowded persons (a, d), inadequate illumination (d, e), specific pose (b), false detection (f) and similar clothes (c, e), which also claim the challenges of person search.

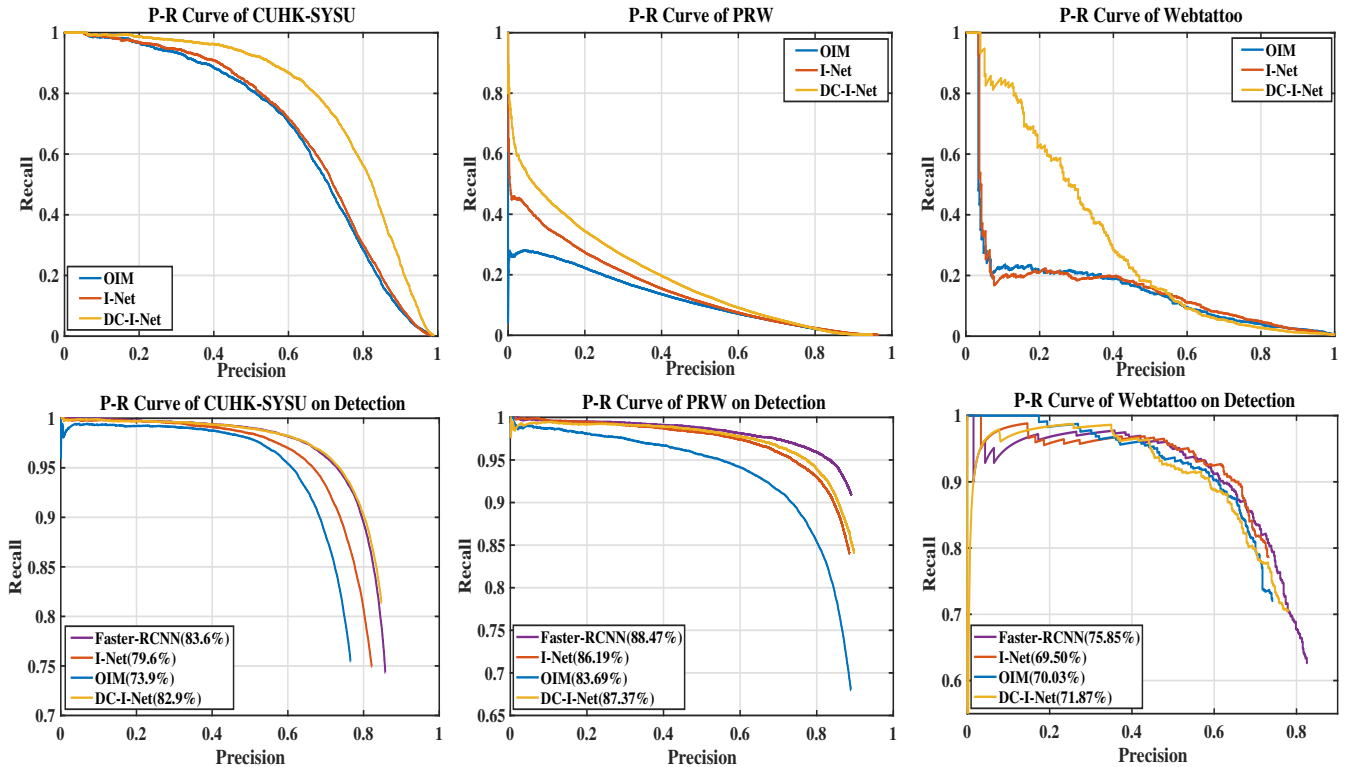


Fig. 8. The Precision-Recall (P-R) curves of different methods on CUHK-SYSU (left), PRW (middle) and Webtattoo (right) datasets for object (i.e., person and tattoo) search and object detection, respectively. The first row shows the P-R curves of object search. The second row presents the P-R curves of detection, in which the AP value of detection is presented in the legend.

Re-ID, the LOMO [40]+XQDA [40], bag of words vector (BOW) [66], IDE_{det} , and CWS [45] are considered. Therefore, 14 methods by combining the separated detection methods and separated Re-ID methods together are compared in this section. Note that in separated detection, for the RCNN based DPM and ACF detectors, AlexNet is implemented as the backbone according to [45].

5.3.2 Experimental Results

The results on the PRW dataset are shown in Table 2, from which we see that our DC-I-Net achieves the best results. Specifically, our method outperforms the SOTA OIM [51] by 10.5% in mAP and 5.2% in top-1 accuracy, which gets similar incremental with the experiments on the CUHK-SYSU [57]. Our method also outperforms the SOTA NPSM [49] by 7.6%

and 2.0% in mAP and top-1 accuracy, respectively. It is worth noting that the end-to-end jointly trained models consistently outperform the separately trained models, which shows the effectiveness of joint multi-task learning for person search. Besides, our DC-I-Net is also much superior to I-Net, which demonstrates the effectiveness of the strategy of divide and conquer of each task in our joint multi-task integration framework. The Precision-Recall (P-R) curves on PRW dataset in person detection and person search are presented in Fig. 8 (middle). We observe that the proposed DC-I-Net shows the best performance for person search task. For detection task, our model is also approaching the state-of-the-art performance of single task of Faster-RCNN detection [11], and outperforms other SOTA models.

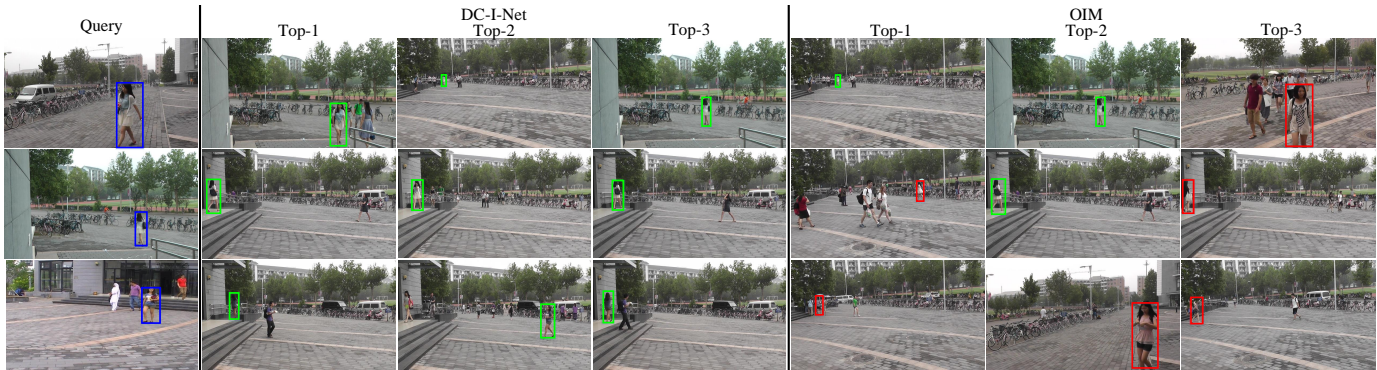


Fig. 9. Retrieval of our DC-I-Net model (middle) and the OIM (right) with the PRW dataset by given three queries (left). The top 3 images with respect to the highest similarity scores of each model are shown. The blue boxes represent the target query person (probe), the green boxes mean correct matches, and the red boxes mean incorrect matches.

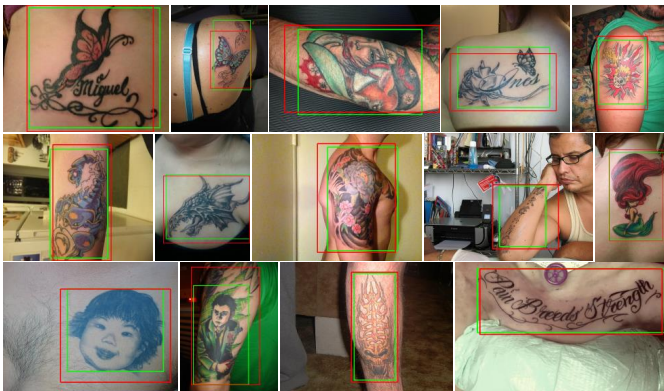


Fig. 10. Detection results of the Webtattoo dataset by using DC-I-Net model. The red boxes represent the predicted bounding boxes and the green boxes represent the ground-truth bounding boxes.

5.3.3 Visualization Results of Person Search

The person search results on PRW dataset are shown in Fig. 9. The pose and resolution variations show the challenge of this dataset. From the top 3 retrieval images with respect to each query, we see that our proposed model can have much better search results than OIM. From the comparison, we see that our model is more robust to both pose and resolution variations.

5.4 Experiments on WebTattoo Dataset

The Webtattoo dataset [58] is proposed search the images containing the same tattoos from the gallery image set as the query (probe) image. Therefore, this dataset is used for object (tattoo) search instead of person search. In this section, we present the results of tattoo detection and tattoo search. Considering that there is few work on tattoo search based on the compared models, for convenience, we compare with the SOTA OIM model [51] because of its open-source advantage. In the experiments, the OIM [51], I-Net, and DC-I-Net are trained and tested without the 300K background tattoo images in the gallery set.

5.4.1 Performance of Tattoo Detection

The detection results based on our DC-I-Net are shown in Fig. 10, from which we could see that the detection of our model approaches the manually annotated ground-truth bounding boxes. Additionally, for comparison of the detection performance, the Precision-recall curves of different models on the Webtattoo dataset for detection are shown in Fig. 8 (right, the second row). We see that the proposed DC-I-Net shows better performance than OIM and I-Net, and it approaches the SOTA Faster-RCNN detection [11], a single task detection model.

5.4.2 Performance of Tattoo Search

To show the importance of the detection part of end-to-end models, the results including mAP, top-1, top-5 and top-10 accuracies with and without (w/o) the detection part are presented in Table 3.

From this table, we observe that the proposed DC-I-Net significantly outperforms I-Net and OIM for both cases. Also, the models with the detection always get better performance than that without the detection part, because the detection enables the model focuses on the tattoo region in the images, such that the performances are dramatically increased sharply with the detection module. With the detection module, our DC-I-Net outperforms OIM and I-Net by 10.3% and 6.8% in mAP, and 11.5% and 7.0% in top-1 accuracy, respectively. Further, the P-R curves of tattoo search on the test set as is shown in Fig. 8 (right, the first row) also demonstrate the clear superiority of the DC-I-Net.

The Tattoo image search results are visualized in in Fig. 11, in which the top 5 images with high similarity scores with respect to each query image are presented. The bounding boxes in the images are the automatically detected object regions of interest for matching and retrieval. Note that in image search task, the whole images rather than the object regions are fed into the model.

6 MODEL ANALYSIS AND DISCUSSION

In the section, to have a deep insight on the effectiveness of the proposed models, the model analysis and discussion are presented based on the CUHK-SYSU [57] dataset.



Fig. 11. Retrieval results of our model with Webtattoo dataset. The top 5 retrieval results for each query are shown. The blue boxes represent the target query person (probe), the green boxes mean correct matches, and the red boxes mean incorrect matches.

TABLE 3
Tattoo search results of different models on Webtattoo dataset

With Detection	mAP(%)	Top-1(%)	Top-5(%)	Top-10(%)
OIM [51]	38.2	39.5	55.5	60.5
I-Net	41.7	44.0	61.0	67.5
DC-I-Net	48.5	51.0	66.5	71.5
w/o Detection	mAP(%)	Top-1(%)	Top-5(%)	Top-10(%)
OIM [51]	21.5	23.0	40.0	41.5
I-Net	23.5	25.5	39.5	44.0
DC-I-Net	30.4	31.0	51.5	64.0

TABLE 4
Performance comparison between I-Net and DC-I-Net based on HEP and C^2 HEP losses, respectively. Note that without special indication, HEP is deployed in I-Net and C^2 HEP is deployed in DC-I-Net.

Loss Type	mAP(%)	Top-1(%)	Top-5(%)	Top-10(%)
I-Net	79.5	81.5	92.2	94.6
I-Net (C^2 HEP)	80.9	83.4	94.1	95.2
DC-I-Net (HEP)	81.0	83.0	93.2	95.6
DC-I-Net	83.7	85.8	94.3	96.1

6.1 Ablation Study of Model Losses

6.1.1 Discussion on the HEP and C^2 HEP Losses

Compared to the I-Net, two key improvements are made in DC-I-Net, including the new network structure and the identity classification loss (i.e., C^2 HEP). We therefore conduct study of the performance gap between HEP in I-Net and C^2 HEP in DC-I-Net. By alternatively changing the softmax guided identity classification loss in each network, the results are presented in Table 4. From the table, we observe that by changing the HEP loss into the C^2 HEP loss in the I-Net, the mAP and top-1 accuracy are improved by 1.4% and 1.9%, respectively. Similarly, by changing the C^2 HEP loss into HEP loss in the DC-I-Net, the mAP and top-1 accuracy are degraded by 2.7% and 2.8%, respectively. These results demonstrate the effectiveness of the newly proposed C^2 HEP loss over HEP loss. Note that the HEP is in default in I-Net and C^2 HEP is in default in DC-I-Net. By comparing the results between I-Net (C^2 HEP) and DC-I-Net, or between I-Net and DC-I-Net (HEP), the increased performance clearly reflects the advantage benefiting from the newly proposed architecture in DC-I-Net.

Remarks. We analyze here the reason why C^2 HEP is effective. The HEP and C^2 HEP losses are supervised by the ground-truth identities. Since person search is an open set problem, we observe the the identity classification accuracy of the training set during iterations. The HEP and the traditional softmax loss are trained for comparing with our C^2 HEP loss. For fair comparison, the ground-truth bounding boxes of only labeled target persons in each image are

TABLE 5
Performance comparisons of different joint losses: metric loss and identity discrimination loss based on the DC-I-Net architecture

Loss Type	mAP(%)	Top-1(%)	Top-5(%)	Top-10(%)
Triplet+HEP	67.8	69.6	87.6	92.2
OLP only	81.3	82.9	93.9	96.0
C^2 HEP only	82.2	84.7	94.0	96.0
OLP + HEP	81.9	83.9	93.9	95.6
OLP + C^2 HEP	83.7	85.8	94.3	96.1

cropped and fed into the model to get the classification accuracy in each iteration. The performance variations based on HEP and C^2 HEP with iterations during the training phase are shown in Fig. 12 (a) and (b), from which we see that the newly proposed C^2 HEP loss significantly outperforms the HEP during the training phase, for both network architectures. We know that in the CUHK-SYSU dataset, 5532 different identities are labeled for training. In each iteration, the input images only contain a very few identities (few-shots) due to the small number of input images, which means that the model training should require thousands of iterations to traverse almost all the identities of the dataset. Therefore, the weights of the old HEP loss is hard to be properly trained, which thus leads to worse performance. In the new C^2 HEP loss, by taken into account the class centers computed via the input features, the identity discrimination can be easily captured because of full-shot property.

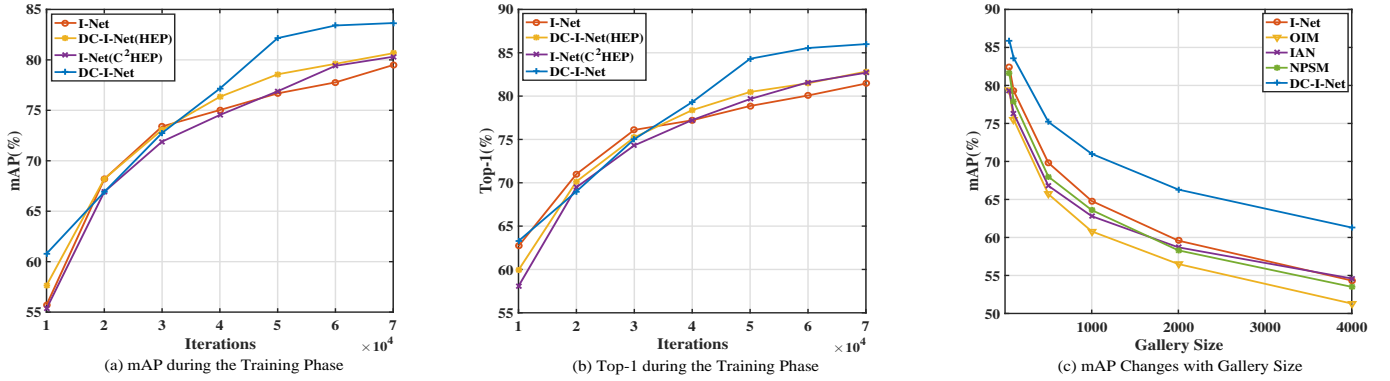


Fig. 12. Discussions on the model. (a) Performance variation (mAP) with different gallery sizes for different models. (b) Performance comparison (mAP) between DC-I-Net and I-Net under different identity classification losses (C²HEP and HEP) with respect to each iteration during the training phase. (c) Performance comparison (top-1 accuracy) between DC-I-Net and I-Net under different identity classification losses (C²HEP and HEP) with respect to each iteration during the training phase. Note that without special indication, I-Net is deployed with HEP loss and DC-I-Net is deployed with C²HEP loss in default.

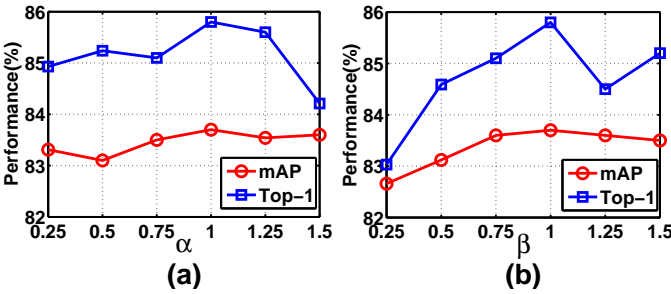


Fig. 13. Sensitivity analysis of the loss weights in Eq. (10). We fix one parameter and adjust another one in experiments. (a) $\beta = 1$; (b) $\alpha = 1$.

6.1.2 Discussion on Joint Metric Loss and Identity Loss

In order to study the effectiveness of the proposed losses, including the metric based OLP loss and the identity based HEP/C²HEP loss. For comparison, we train our network based on different combinations of the metric loss (triplet loss, OLP) and identity loss (HEP, C²HEP). Note that the network structure is the same as the DC-I-Net. The results are presented in Table 5.

From the table, we see that the traditional triplet loss together with the softmax guided cross-entropy loss achieves much worse performance in mAP and top-1 accuracy than ours. With the proposed OLP loss only or C²HEP loss only, much better performances can be achieved. The worse performance of triplet loss is mainly caused by the stagnation problem during training the end-to-end person search model. Therefore, the proposed softmax guided OLP loss can well overcome the stagnation problem of the triplet loss. With the proposed C²HEP based identity loss only, the performance is comparable to OLP based metric loss only. By combining them together, the proposed DC-I-Net achieves the best performance, 83.7% and 85.8% in mAP and top-1 accuracy, respectively. The effectiveness of the proposed losses are verified in this section.

6.1.3 The Analysis of Loss Weight

In this section, we adjust the weights of OLP loss and C²HEP loss of DC-I-Net during the training phase. We fix

TABLE 6
Performance with different feature dictionary size stored in OLP based on the DC-I-Net architecture.

mAP (%)	20 × 128	40 × 128	60 × 128
OLP only	81.4	81.3	81.8
OLP+C ² HEP	82.5	83.7	83.1
Top-1 (%)	20 × 128	40 × 128	60 × 128
OLP only	83.5	82.9	83.4
OLP+C ² HEP	84.7	85.8	85.2

the α as 1.0 in the Eq. (10) and adjust the β from 0.25 to 1.5 to study the influence and vice versa. The backbone of DC-I-Net is set as VGG16 in the experiments. The mAP and Top-1 with respect to different loss weights are shown in the Fig. 13. From the results, we observe that the trade-off parameters have slight impact and the best performance can be easily achieved when $\alpha = 1$ and $\beta = 1$.

6.2 Analysis of the Feature Dictionary in OLP and Priority Classes in C²HEP

6.2.1 Impact of the Feature Dictionary

In the proposed OLP loss, a feature dictionary is deployed to generate negative pairs, which may preferably restrict the positive pairs for effective metric learning. Therefore, the impact of the dictionary size, i.e. the number of features stored, is studied. The model is trained either with only OLP or with both OLP and C²HEP loss.

The number of features stored in the feature dictionary of OLP depends on the minibatch size (i.e., 128). In this study, different times (i.e., 20, 40 and 60) the minibatch size are determined as the dictionary size (i.e., 20×128, 40×128, 60×128) and tested, respectively. Note that 40×128 means that the number of features stored in the dictionary is 40 times the minibatch size (128), i.e. 5120. The analysis results are presented in Table 6, in which we see that the impact of the dictionary size is slight. With both the OLP and C²HEP loss functions implemented, i.e., DC-I-Net, the model gets the best results when the dictionary size is set as 40×128.

From Table 6, two perspectives can be observed. 1) From the case of OLP loss only trained model, increasing the

TABLE 7
Performance comparisons by using different numbers of selected priority classes based on the DC-I-Net architecture.

	mAP(%)	Top-1(%)	Top-5(%)	Top-10(%)
50	82.5	84.8	94.1	96.0
100	83.7	85.8	94.3	96.1
1000	83.1	85.0	94.6	95.8
5532	82.9	83.7	93.8	95.6

feature dictionary size does not improve the performance, which is because the OLP loss is hard to contain all identities in the dataset and even the number of stored features grows, the number of identities in the feature dictionary is not significantly increased. 2) From the case of both OLP and C²HEP losses trained model, the best performance with a suitable size of feature dictionary can be achieved, which is because the C²HEP loss can explore all the labeled identities in the dataset. As a result, the dictionary size is set as 40×128 in experiments.

6.2.2 Impact of the Number of Priority Classes

The essence of C²HEP lies in the class center guided hard example priority mechanism. Therefore, utilizing different number of priority classes for computing the C²HEP loss is studied in experiments. Note that, in this study, both OLP and C²HEP loss functions are used for implementation. Specifically, in the experiments, the number of the selected priority classes is set as 50, 100, 1000, and 5532, respectively. Note that the number 5532 means that all the labeled persons are used in loss computation without consideration of the priority classes. The results are presented in Table 7, from which we observe that the model trained with 100 selected priority classes outperforms the model without the hard example priority strategy (i.e., 5532) by 0.8% in mAP and 2.1% in top-1 accuracy. The effectiveness of the proposed hard example priority strategy is verified. Additionally, with the increasing of the number of priority classes, the performance is slightly degraded which is due to that the attention on the really hard classes may be reduced. In contrast, if the number of priority classes is too small, the model may not well explore the identities of the whole dataset such that the performance is not good. Therefore, a suitable number of priority classes is required, which is set as 100 in experiments.

6.3 Analysis of the Retrieval Performance

6.3.1 Impact of Gallery Size

The person search is essentially an image retrieval task, which therefore becomes more challenging especially when the size of gallery set (retrieval pool) increases. This section presents a study on the impact of different gallery size. Specifically, in the experiments, we vary the gallery size as 50, 100, 500, 1000, 2000, and 4000, respectively, and test the mAP of different models including DC-I-Net, I-Net, the OIM [51], NPSM [49], and IAN [67] for each gallery size. The retrieval performance variation curves are shown in Fig. 12 (c), from which we can see that with the increasing gallery size of the test set, the performances of all methods

are degraded. It is worthy noting that our proposed DC-I-Net has a relatively slower degradation speed, but always outperforms other methods with respect to each gallery size. In fact, it is common that the difficulty in finding the query

TABLE 8
Performance analysis with different numbers of input images at every iteration based on different losses.

	OLP+C ² HEP		Contrastive Loss	
	mAP(%)	Top-1(%)	mAP(%)	Top-1(%)
2-input	86.2	86.5	48.7	45.0
4-input	85.9	87.0	54.4	54.7
8-input	85.8	85.9	60.4	60.7

person is growing with the increasing gallery size because of the enhanced inter-similarity.

6.3.2 Performance Variation With Training Iterations

This section presents the retrieval performance variation of mAP and top-1 accuracy by using the proposed DC-I-Net and the previous I-Net, both of which are trained for 70K iterations. The retrieval performance variation of our models during the training phase is shown in Fig. 12 (a) and (b), where both the mAP and top-1 accuracy increase with the training iterations. DC-I-Net shows faster upward trend than I-Net. Note that the learning rate is reduced at the 40K iteration in the training phase.

6.4 Analysis of Different Numbers of Input Images

The insufficient number of input images for each iteration leads to the stagnation problem for the traditional loss function (e.g. contrastive loss). In this part, we change the number of input images of each iteration for ablation analysis. The resnet50 based DC-I-Net architecture is implemented based on the proposed losses and the contrastive loss. Specifically, we set the number of input training images as 2, 4, and 8, respectively in each iteration in this experiment. The results are shown in Table 8, from which we can clearly see that when the number of input images is small, the traditional contrastive loss even does not work. The reason is that the model based on contrastive loss encounters the training stagnation problem when the number of input images is set as 2. As the number of input images increases, the stagnation problem is alleviated and the performance is progressively improved.

On the contrary, our proposed losses always work well. The proposed OLP and C²HEP loss functions have well reduced the influence of the stagnation problem, so that the number of input images does not have much impact on our model. This fully verifies the motivation and effectiveness of the proposed OLP and C²HEP losses.

7 CONCLUSION

In this paper, we propose an end-to-end tasks-integrated network (I-Net) for user-friendly image search, by jointly modeling the detection and retrieval tasks in a unified framework. While I-Net does not consider the task specification and the inherent problem of small number of input images in training, we have further proposed an essentially

improved integrated network with the philosophy of divide and conquer, called DC-I-Net. Two merits are naturally formulated: 1) the task specification can be explored and more accurate object proposals are used to effectively train the re-identifier. 2) A novel class center guided hard example priority (C^2 HEP) loss is proposed by utilizing the class centers computed via the timely updated input features, which well overcomes the inherent few-shot (i.e., very few identities) problem during training of image search model. The proposed models outperform the state-of-the-art task-integrated and task-separated image search models on three widely used benchmark datasets, such as CUHK-SYSU [57], PRW [45] and Webtattoo [58].

REFERENCES

- [1] X. Wang, "Intelligent multi-camera video surveillance: A review," *Pattern Recognition Letters*, vol. 34, no. 1, pp. 3–19, 2013.
- [2] B. Song, A. T. Kamal, C. Soto, C. Ding, J. A. Farrell, and A. K. Roychowdhury, "Tracking and activity recognition through consensus in distributed camera networks," *IEEE Trans. Image Processing*, vol. 19, no. 10, pp. 2564–2579, 2010.
- [3] S. Z. Chen, C. C. Guo, and J. H. Lai, "Deep ranking for person re-identification via joint representation learning," *IEEE Trans. Image Processing*, vol. 25, no. 5, pp. 2353–2367, 2016.
- [4] C. Cao, Y. Wang, J. Kato, G. Zhang, and K. Mase, "Solving occlusion problem in pedestrian detection by constructing discriminative part layers," in *WACV*, 2017, pp. 91–99.
- [5] H. Song, W. Wang, J. Wang, and R. Wang, "Collaborative deep networks for pedestrian detection," in *IEEE Third International Conference on Multimedia Big Data*, 2017, pp. 146–153.
- [6] W. Ouyang, X. Zeng, and X. Wang, "Modeling mutual visibility relationship in pedestrian detection," in *CVPR*, 2013, pp. 3222–3229.
- [7] Y. Shen, R. Ji, C. Wang, X. Li, and X. Li, "Weakly supervised object detection via object-specific pixel gradient," *IEEE Trans. Neural Networks and Learning Systems*, vol. 29, no. 12, pp. 5960–5970, 2018.
- [8] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *ICCV*, 2017, pp. 2961–2969.
- [9] W. Li, R. Zhao, and X. Wang, "Human reidentification with transferred metric learning," in *ACCV*, 2012, pp. 31–44.
- [10] W. Li and X. Wang, "Locally aligned feature transforms across views," in *CVPR*, 2013, pp. 3594–3601.
- [11] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2017.
- [12] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *CVPR*, 2015, pp. 815–823.
- [13] D. Cheng, Y. Gong, S. Zhou, J. Wang, and N. Zheng, "Person re-identification by multi-channel parts-based cnn with improved triplet loss function," in *CVPR*, 2016, pp. 1335–1344.
- [14] J. Liu, Z. J. Zha, Q. I. Tian, D. Liu, T. Yao, Q. Ling, and T. Mei, "Multi-scale triplet cnn for person re-identification," in *ACM MM*, 2016, pp. 192–196.
- [15] S. Ding, L. Lin, G. Wang, and H. Chao, "Deep feature learning with relative distance comparison for person re-identification," *Pattern Recognition*, vol. 48, no. 10, pp. 2993–3003, 2015.
- [16] C. Song, Y. Huang, W. Ouyang, and L. Wang, "Mask-guided contrastive attention model for person re-identification," in *CVPR*, 2018, pp. 1179–1188.
- [17] Z. He and L. Zhang, "End-to-end detection and re-identification integrated net for person search," in *ACCV*, 2018, pp. 349–364.
- [18] P. Dollár, S. Belongie, S. Belongie, and P. Perona, "Fast feature pyramids for object detection," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 36, no. 8, p. 1532, 2014.
- [19] W. Nam, P. Dollar, and J. H. Han, "Local decorrelation for improved pedestrian detection," in *NIPS*, 2014, pp. 1–9.
- [20] S. Zhang, R. Benenson, and B. Schiele, "Filtered channel features for pedestrian detection," in *CVPR*, 2015, pp. 1751–1760.
- [21] P. Dollar, Z. Tu, P. Perona, and S. Belongie, "Integral channel features," in *British Machine Vision Conference*, 2009.
- [22] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *CVPR*, 2005.
- [23] Y. Tian, P. Luo, X. Wang, and X. Tang, "Pedestrian detection aided by deep learning semantic tasks," in *CVPR*, 2014, pp. 5079–5087.
- [24] L. Zhang, L. Lin, X. Liang, and K. He, "Is faster r-cnn doing well for pedestrian detection?" in *ECCV*, 2016, pp. 443–457.
- [25] J. Xu, R. Zhao, F. Zhu, H. Wang, and W. Ouyang, "Attention-aware compositional network for person re-identification," in *CVPR*, 2018, pp. 2119–2128.
- [26] Z. Zheng, L. Zheng, and Y. Yang, "Pedestrian alignment network for large-scale person re-identification," *IEEE Trans. Circuits and Systems for Video Technology*, 2018.
- [27] Z. Zhong, L. Zheng, D. Cao, and S. Li, "Re-ranking person re-identification with k-reciprocal encoding," in *CVPR*, 2017.
- [28] T. Xiao, H. Li, W. Ouyang, and X. Wang, "Learning deep feature representations with domain guided dropout for person re-identification," in *CVPR*, 2016.
- [29] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *ICCV*, 2017, pp. 3754–3762.
- [30] H.-X. Yu, A. Wu, and W.-S. Zheng, "Cross-view asymmetric metric learning for unsupervised person re-identification," in *ICCV*, 2017, pp. 994–1002.
- [31] W.-S. Zheng, X. Li, T. Xiang, S. Liao, J. Lai, and S. Gong, "Partial person re-identification," in *ICCV*, 2015, pp. 4678–4686.
- [32] A. J. Ma, P. C. Yuen, and J. Li, "Domain transfer support vector ranking for person re-identification without target camera label information," in *ICCV*, 2013, pp. 3567–3574.
- [33] W. Deng, L. Zheng, Q. Ye, G. Kang, Y. Yang, and J. Jiao, "Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification," in *CVPR*, 2017.
- [34] F. Liu and L. Zhang, "View confusion feature learning for person re-identification," in *ICCV*, 2019.
- [35] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in *CVPR*, 2015, pp. 3908–3916.
- [36] Z. Zheng, L. Zheng, and Y. Yang, "A discriminatively learned cnn embedding for person re-identification," *ACM Trans. Multimedia Computing, Communications, and Applications*, vol. 14, no. 1, 2018.
- [37] W. Chen, X. Chen, J. Zhang, and K. Huang, "A multi-task deep network for person re-identification," in *AAAI*, 2017.
- [38] X. Wang, G. Doretto, T. Sebastian, J. Rittscher, and P. H. Tu, "Shape and appearance context modeling," in *ICCV*, 2007.
- [39] D. Gray and T. Hai, "Viewpoint invariant pedestrian recognition with an ensemble of localized features," in *ECCV*, 2008.
- [40] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *CVPR*, 2015, pp. 2197–2206.
- [41] P. M. Roth, P. Wohlhart, M. Hirzer, M. Kostinger, and H. Bischof, "Large scale metric learning from equivalence constraints," in *CVPR*, 2012, pp. 2288–2295.
- [42] X. Li, W. S. Zheng, X. Wang, and T. Xiang, "Multi-scale learning for low-resolution person re-identification," in *ICCV*, 2015, pp. 3765–3773.
- [43] Y. Sun, L. Zheng, Y. Li, Y. Yang, Q. Tian, and S. Wang, "Learning part-based convolutional features for person re-identification," *IEEE Trans. Pattern Analysis and Machine Intelligence*, 2019.
- [44] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: a deep quadruplet network for person re-identification," in *CVPR*, 2017, pp. 403–412.
- [45] L. Zheng, H. Zhang, S. Sun, M. Chandraker, Y. Yang, and Q. Tian, "Person re-identification in the wild," in *CVPR*, 2016, pp. 1367–1376.
- [46] X. Lan, X. Zhu, and S. Gong, "Person search by multi-scale matching," in *ECCV*, 2018, pp. 536–552.
- [47] D. Chen, S. Zhang, W. Ouyang, J. Yang, and Y. Tai, "Person search via a mask-guided two-stream cnn model," in *ECCV*, 2018, pp. 734–750.
- [48] C. Han, J. Ye, Y. Zhong, X. Tan, C. Zhang, C. Gao, and N. Sang, "Re-id driven localization refinement for person search," in *ICCV*, 2019, pp. 9814–9823.
- [49] H. Liu, J. Feng, Z. Jie, K. Jayashree, B. Zhao, M. Qi, J. Jiang, and S. Yan, "Neural person search machines," in *ICCV*, 2017, pp. 493–501.

- [50] Y. Yan, Q. Zhang, B. Ni, W. Zhang, M. Xu, and X. Yang, "Learning context graph for person search," in *CVPR*, 2019, pp. 2158–2167.
- [51] T. Xiao, S. Li, B. Wang, L. Lin, and X. Wang, "Joint detection and identification feature learning for person search," in *CVPR*, 2017.
- [52] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv*, 2014.
- [53] R. Girshick, "Fast r-cnn," in *ICCV*, 2015, pp. 1440–1448.
- [54] K. Sohn, "Improved deep metric learning with multi-class n-pair loss objective," in *NIPS*, 2016.
- [55] Y. Wen, K. Zhang, Z. Li, and Y. Qiao, "A discriminative feature learning approach for deep face recognition," in *ECCV*, 2016.
- [56] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *NIPS*, 2017.
- [57] T. Xiao, S. Li, B. Wang, L. Lin, and X. Wang, "End-to-end deep learning for person search," *arXiv*, 2016.
- [58] H. Han, J. Li, A. K. Jain, S. Shan, and X. Chen, "Tattoo image search at scale: Joint detection and compact representation learning," *IEEE Trans. Pattern Analysis and Machine Intelligence*, 2019.
- [59] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *ACM MM*, 2014, pp. 675–678.
- [60] M. Ngan, G. Quinn, and P. Grother, "Tattoo recognition technology-challenge (tatt-c): Outcomes and recommendations," in *Tech. Rep.*, 2018.
- [61] Q. Xu, S. Ghosh, X. Xu, Y. Huang, and A. Kong, "Tattoo detection based on cnn and remarks on the nist database," in *ICB*, 2016.
- [62] T. Hrkač, K. Brkić, and Z. Kalafatić, "Tattoo detection for soft biometric de-identification based on convolutional neural networks," in *The 1st OAGM-ARW Joint Workshop-Vision Meets Robotics*, 2016.
- [63] B. Yang, J. Yan, Z. Lei, and S. Z. Li, "Convolutional channel features," in *ICCV*, 2015, pp. 82–90.
- [64] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016, pp. 770–778.
- [65] R. Zhao, W. Ouyang, and X. Wang, "Unsupervised saliency learning for person re-identification," in *CVPR*, 2013, pp. 3586–3593.
- [66] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *ICCV*, 2015, pp. 1116–1124.
- [67] J. Xiao, Y. Xie, T. Tillo, K. Huang, Y. Wei, and J. Feng, "Ian: The individual aggregation network for person search," *Pattern Recognition*, vol. 87, pp. 332–340, 2017.
- [68] X. Chang, P.-Y. Huang, Y.-D. Shen, X. Liang, Y. Yang, and A. G. Hauptmann, "Rcaa: Relational context-aware agents for person search," in *ECCV*, 2018, pp. 84–100.
- [69] P. F. Felzenszwalb, R. B. Girshick, D. Mcallester, and D. Ramanan, "Object detection with discriminatively trained part-based models." *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 32, no. 9, pp. 1627–1645, 2010.