# Deep High-Resolution Representation Learning for Cross-Resolution Person Re-identification

Guoqing Zhang, *Member, IEEE,* Yu Ge, Zhicheng Dong, Hao Wang, Yuhui Zheng, and Shengyong Chen *Senior Member, IEEE*

*Abstract*—**Person re-identification (re-ID) tackles the problem of matching person images with the same identity from different cameras. In practical applications, due to the differences in camera performance and distance between cameras and persons of interest, captured person images usually have various resolutions. We name this problem as Cross-Resolution Person Re-identification which brings a great challenge for matching correctly. In this paper, we propose a Deep High-Resolution Pseudo-Siamese Framework (PS-HRNet) to solve the above problem. Specifically, in order to restore the resolution of low-resolution images and make reasonable use of different channel information of feature maps, we introduce and innovate VDSR module with channel attention (CA) mechanism, named as VDSR-CA. Then we reform the HRNet by designing a novel representation head to extract discriminating features, named as HRNet-ReID. In addition, a pseudo-siamese framework is constructed to reduce the difference of feature distributions between low-resolution images and high-resolution images. The experimental results on five cross-resolution person datasets verify the effectiveness of our proposed approach. Compared with the state-of-the-art methods, our proposed PS-HRNet improves 3.4%, 6.2%, 2.5%,1.1% and 4.2% at Rank-1 on MLR-Market-1501, MLR-CUHK03, MLR-VIPeR, MLR-DukeMTMC-reID, and CAVIAR datasets, respectively. Our code is available at https://github.com/zhguoqing.**

*Index Terms*—**Cross-resolution person re-identification, super-resolution, high-resolution network, pseudo-siamese framework, deep learning.**

## I. INTRODUCTION

**P**ERSON re-identification (re-ID) intends to match person images with the same identity across images captured by various cameras. Re-ID has become the spotlight in the field of machine learning and computer vision owing to its wide practicability in recent years [1], [2], [3], [4], [5], [6], [7], [8], [9]. Driven by recent advances of deep learning, existing researches of re-ID focus on designing deep feature extraction networks to improve the matching accuracy of re-ID [10], [11], [12], [13]. Although these approaches have achieved satisfactory performance and alleviated the influence of person pose changes, background clutters or part occlusions to a certain extent, these methods are usually on the basis of

G. Zhang is with the School of Computer and Software, Nanjing University of Information Science and Technology, Nanjing, 210044, China, and also with the Engineering Research Center of Digital Forensics, Ministry of Education, Nanjing University of Information Science and Technology (e-mail: xiayang14551@163.com). (Corresponding author: Yuhui Zheng).

Y. Ge, Z. Dong, H. Wang, and Y. Zheng are with the School of Computer and Software, Nanjing University of Information Science and Technology, China, 210044 (e-mail: gy1328447669@gmail.com, dzc2000919@gmail.com, btnode3@gmail.com, zheng_yuhui@nuist.edu.cn).

S. Chen is with the School of Computer Science and Engineering, Tianjin University of Technology, Tianjin 300384, China (e-mail: sy@ieee.org).



Fig. 1. Illustration of the difference between (a) traditional person re-ID task in the ideal scenario and (b) cross-resolution person re-ID task. Compared with high-resolution (HR) query images and gallery images, low-resolution (LR) query images contain less fine-grained details, which causes a significant reduction in recognition accuracy and brings great challenges to the work of matching.

the prerequisite that the gallery images and the query images possess the same resolution and sufficient fine-grained details.

However, such prerequisite is difficult to guarantee in practical applications. The problem of matching the same person images with different resolutions is named as Cross-Resolution Person Re-identification [13], [14], [15], [16].

Fig. 1 shows the difference between (a) the person re-ID task in an ideal condition and (b) the cross-resolution person re-ID task. Ideally, query images maintain the same high-resolution (HR) as gallery images. However, due to the differences in camera performance and distance between probes and target pedestrians, captured query images often possess lower resolution than gallery images. The lack of image information makes the traditional re-ID methods incapable of effectively extracting the discriminant features of images for matching, which has become a stumbling block to the development of re-ID task.

In order to settle the above problem, many cross-resolution person re-ID algorithms have been put forward in recent years. In early work, the main idea plans to explore the common feature representation space of HR and LR images by using metric learning or dictionary learning methods, such as [14], [13], [17], [18]. However, the performance of these methods are restricted due to the incapability of recovering the information lost in LR images. Later, some researchers try to introduce super-resolution (SR) technology into cross-resolution person re-ID. SING [16] first applies SRCNN [19] as the resolution recovery module and jointly trains the SRCNN sub-network and the re-ID sub-network. Since then, different SR networks, such as SRGAN [20] and FFSR [21], are introduced as the

resolution recovery module to further optimize the framework. Recently, some new methods represented by INTACT [40] have been proposed, and more novel and effective mechanisms have been applied to raise the detection accuracy to a new level. These methods have achieved significant performance improvement, but it is still far below the practical application standard.

Through detailed comparison and analysis of numerous recent cross-resolution person re-ID methods, we gradually discovered some commonalities contained in them. The most conspicuous is that the existing approaches almost all use convolutional neural networks with the property of down-sampling such as ResNet [53] as feature extraction networks. We believe that this is the most detrimental factor in the existing methods. Using such networks as the feature extraction backbone will inevitably cause further loss of fine-grained information from low-resolution images. Besides, excessive emphasis on low-resolution image reconstruction seems to have formed a stereotyped thinking pattern. In fact, through experiments, we found that complex super-resolution networks may not perform better than a simple one in cross-resolution person re-ID task under certain circumstances. More energy should be devoted to the study of deep semantic information and feature information extracted from low-resolution images.

In this paper, we propose the PS-HRNet to solve the limitations analyzed above. Firstly, we further improve the super-resolution capability of VDSR [22] in terms of deep semantic information learning by adding the channel attention mechanism, and name the modified SR module as VDSR-CA. Besides, based on the finding that the unique parallel architecture of HRNet [23] is helpful to alleviate the impact of resolution difference, we utilize the HRNet as the feature extraction network. Here we propose the HRNet-ReID to capture multi-resolution features of person images by introducing a novel representation head to HRNet, which can adapt HRNet to the person re-ID problem. In addition, our PS-HRNet adopts a pseudo-siamese framework [65], [66] so as to further decrease the distribution difference between LR image features and HR image features. The training strategy of the whole network is divided into two phases. In the first phase, only the HRNet-ReID module in HR branch of pseudo-siamese framework is trained on traditional HR person re-ID datasets. In the second phase, we use joint training strategy to train both the VDSR-CA module and two HRNet-ReID modules simultaneously on cross-resolution person re-ID datasets.

The outstanding contributions of our work are summarized in the following three points:

• We put forward a feature extraction network named as HRNet-ReID, which combines native HRNet-W32 backbone with a novel representation head designed by us to adapt HRNet to the specific person re-ID mission, overcoming the flaw caused by conventional feature extraction networks in existing methods.

• We construct a pseudo-siamese framework named as PS-HRNet which combines our proposed VDSR-CA and HRNet-ReID to further explore the feature space at a deeper level and successfully reduce the distribution difference between LR image features and HR image features, providing an original solution to the cross-resolution person re-ID problem.

• We have carried out extensive experiments on five cross-resolution person re-ID datasets, and all of them have achieved the highest level in the industry. Compared with the state-of-the-arts, our proposed PS-HRNet improves 3.4%, 6.2%, 2.5%, 1.1% and 4.2% at Rank-1 on MLR-Market-1501, MLR-CUHK03, MLR-VIPeR, MLR-DukeMTMC-reID, and CAVIAR datasets, respectively.

## II. RELATED WORK

In this section, we will roughly introduce the related works concerned with traditional person re-ID and cross-resolution person re-ID, and revisit two utilized core modules and effective pseudo-siamese framework.

### A. Person re-ID

Person re-ID has achieved rapid development in the past decades. A series of methods have been proposed to extract more robust and discriminating feature representations, and overcome the difficulties brought by person pose changes, background clutters or part occlusions. Specifically, to solve pose changes, Liu et al. [26] design a pose-transferable GAN which aims to produce person images with multiple poses for data enhancement. To address background clutter, some methods based on attention mechanism or semantic parsing are proposed. Li et al. [27] apply spatial and channel attention to make network focus on more informative parts. Kalayeh et al. [28] adopt semantic parsing to segment the foreground information and background information to reduce the interference of background. Besides, extensive methods have achieved great progress in occluded re-ID [29], [30], unsupervised re-ID [31], [32], cross-modality re-ID [33], [34], and so on. However, most of existing methods neglect the resolution mismatch problem which is a common situation in practical scenarios.

### B. Cross-Resolution Person re-ID

To solve the resolution inconsistency problem, a few methods have been proposed recently. Previous traditional methods [35], [36] mainly focus on dictionary learning and metric learning, which achieve limited performance due to the lack of detail features in LR images. Encouraged by the flourishing development of convolutional neural networks (CNN) [37] and super-resolution (SR) technology, some SR-based methods are proposed and greatly improve the matching accuracy. For instance, Jiao et al. [16] make the first attempt to combine the SRCNN and re-ID network into one framework, and propose a jointly training strategy. Besides, some methods adopt GANs to further improve the framework. Specifically, Wang et al. [38] adopt SR-GAN repeatedly to build a cascaded structure. Li et al. [39] restore image resolutions and learn the resolution-invariant representations. Recently, Cheng et al. [40] optimize SR-reID joint framework from the perspective of training strategy and achieve the best performance, which enhances the compatibility between two sub-networks by utilizing the underlying association knowledge between SR and re-ID.
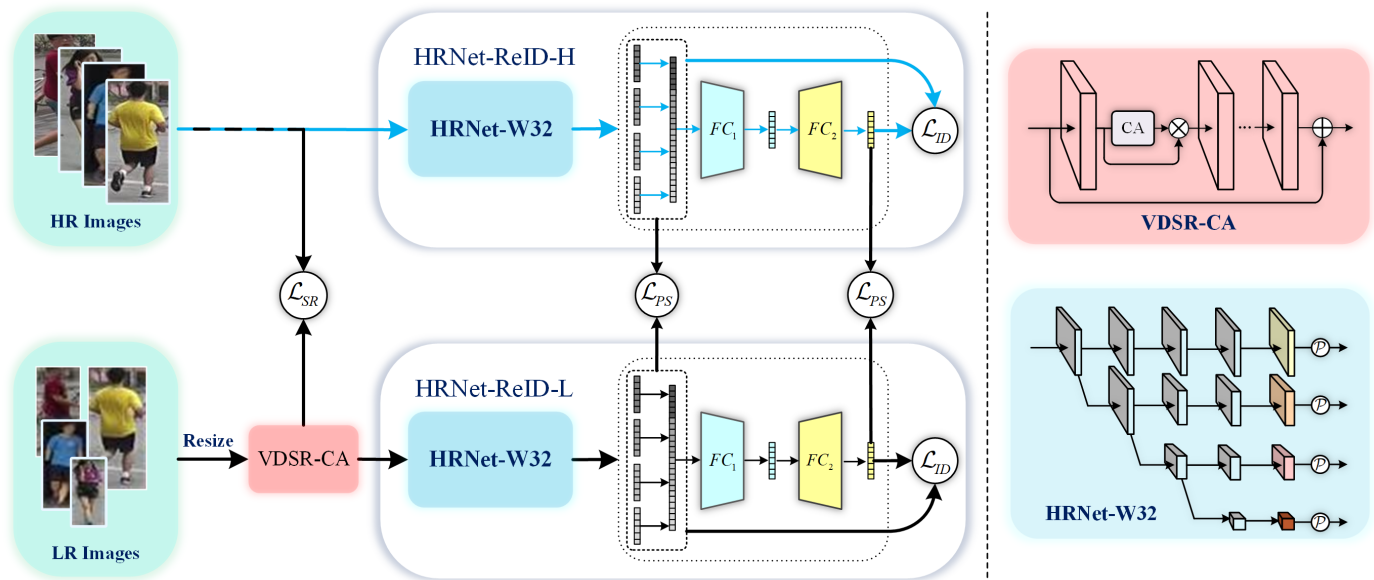
Fig. 2. The unified architecture of Deep High-Resolution Pseudo-Siamese Framework (PS-HRNet) proposed by us for cross-resolution person re-ID task, which is a pseudo-siamese framework consisting of double HRNet-ReID feature extraction networks and one VDSR-CA module. The VDSR-CA is utilized to restore the resolution of input LR images. The HRNet-ReID is designed to adapt the HRNet to the person re-ID task, and extracts discriminating features from restored images. The right side gives a brief sketch of the two types of modules. $\otimes$ and $\oplus$ denote element-wise product and element-wise add, respectively. As the main structure in PS-HRNet, the pseudo-siamese framework is adopted to close the feature distributions of LR and HR images. Our PS-HRNet is trained alternatively in two phases: **(1)** Update the single HRNet-ReID-H with the loss $\mathcal{L}_{ID}$ (Eq. (10)); **(2)** Update the joint multi-task learning loss with the loss $\mathcal{L}_{TOTAL}$ (Eq. (12)). The above two phases are marked with blue and **black** arrows, respectively. (best viewed in color).

Most existing deep learning methods based on SR attach their importance to reconstruct SR images and make the generated images visually closer to the original HR images. However, these methods ignore the distribution differences between the LR image features and HR image features extracted by the feature extraction network separately.

### C. Revisit VDSR and HRNet

As a high-performance super-resolution method, VDSR [22] applies a deeper network to gain in-depth image information and further ameliorates the structure of SRCNN [19]. Motivated by the prevalence and development of residual network (ResNet) , VDSR adopts the residual connection to overcome the difficulty in convergence of deep networks. Therefore, the capacity of VDSR on image reconstruction surpasses SRCNN strikingly.

HRNet [23] is first proposed to deal with the Human Pose Estimation task, and then surpasses all predecessors in other fields such as key point detection, pose estimation and multi-person pose estimation [41], [42]. The core structure of HRNet contains four parallel streams, which is logically presented as follow:

$$\begin{aligned}
\mathcal{N}_{11} \rightarrow \mathcal{N}_{21} &\rightarrow \mathcal{N}_{31} \rightarrow \mathcal{N}_{41} \\
\searrow \mathcal{N}_{22} &\rightarrow \mathcal{N}_{32} \rightarrow \mathcal{N}_{42} \\
&\searrow \mathcal{N}_{33} \rightarrow \mathcal{N}_{43} \\
&\qquad \searrow \mathcal{N}_{44},
\end{aligned} \qquad (1)$$

where $\mathcal{N}_{sr}$ is a sub-stream in the $s$-th phase and $r$ is a resolution index. In order to maintain the high resolution of the same branch while obtaining information from other branches, layers at the junction of different $\mathcal{N}$ in the same $s$-th can propagate all of its information to every sub-stream $\mathcal{N}$ in the

next $s$-th phase. Such unique parallel-cross structure enables the module to acquire and transmit information between networks of different scales. Besides, the structure has the ability to maintain high-resolution feature map for each branch and each phase.

### D. Pseudo-Siamese Framework

The siamese neural network architecture was first proposed to verify the signature of a check at the end of the last century, and achieved satisfactory expectations [24]. It feeds inputs into two identical neural networks that share weights with each other to map the inputs to a brand new feature space and then measure the difference between the inputs under the new feature representation [25].

Inspired by the success of siamese neural network architecture, the pseudo-siamese framework is proposed for various detection and recognition tasks [66]. Different from the weights sharing mechanism of the siamese neural network architecture, the network contained in the latter does not need to share weights, so it can be composed of two identical or different sub-networks, which makes the pseudo-siamese framework possess higher degrees of flexibility and a wider range of application scenarios. The proposal of the pseudo-siamese network brings a novel thinking to traditional classification and comparison tasks.

On the basis of the above existing works, in our method, we add the channel attention (CA) mechanism [43] to VDSR, which makes the network perceive the more informative channels in the process of image reconstruction. In addition, we design a brand-new representation head of HRNet to make full

utilization of person image features. Finally, we structure the pseudo-siamese framework, composed of two different sub-networks, by which high-resolution and low-resolution pedestrian images received respectively to explore the similarity and feature distribution of cross-resolution images in new feature representation.

## III. PROPOSED METHOD

In this section, the proposed PS-HRNet is introduced for the cross-resolution person re-ID problem by giving an overview of its unified architecture first, followed by the details of its main components and the pseudo-siamese framework.

### A. Framework Overview

As illustrated in Fig. 2, our proposed PS-HRNet adopts a pseudo-siamese framework as the global structure, and contains two main modules, i.e., HRNet-ReID and VDSR-CA. In order to clearly clarify the following formulas, we first define the notations of datasets which will be used in this paper. In the training phase, we define $N$ high-resolution (HR) images with associated labels as $\mathcal{D}_h = \{x_h^i, y^i\}_{i=1}^N$, where $x_h^i \in \mathbb{R}^{H \times W \times 3}$. We down-sample each HR image with the down-sample rate $r \in \{2, 3, 4\}$ (i.e. the spatial size of a down-sampled image becomes $\frac{H}{r} \times \frac{W}{r}$). The generated corresponding low-resolution (LR) images are denoted as $\mathcal{D}_l = \{x_l^i, y^i\}_{i=1}^N$.

Our proposed PS-HRNet method has two main objectives. Just as most SR-based methods, one objective is to use the super-resolution reconstruction module to restore the missing detail information in LR images and reduce the visual difference between LR and HR images. In our method, the original VDSR [22] is improved by adding a channel attention block and named as VDSR-CA, which is applied to generate a HR version for each LR image. The other is to minimize the discrepancies in feature distribution between LR and HR images and enhance the matching accuracy by constructing a pseudo-siamese framework. For each pair of input images, we utilize HRNet-ReID as the feature extraction network for each branch to learn HR and LR feature representations, and use losses to promote the reduction in distribution differences of these features. These modules will be illustrated meticulously in following subsections.

### B. VDSR with Channel Attention

In most traditional CNN-based super-resolution modules, each channel of the feature map is treated equally in the process of information transmission between every two feature maps. However, the reality is that the image features contained in different channels of the feature map are various, and these diversities contribute to the recovery of high-frequency features in super-resolution task to a different extent. So it is very meaningful to assign different weights to the channels of feature map to embody the difference between the channels. Therefore, we adopt a channel attention (CA) mechanism proposed in RCAN [43] to redistribute the characteristics of different channels of the feature map to enhance the recovery ability of SR networks. The formulaic representation of the CA mechanism is as follows:

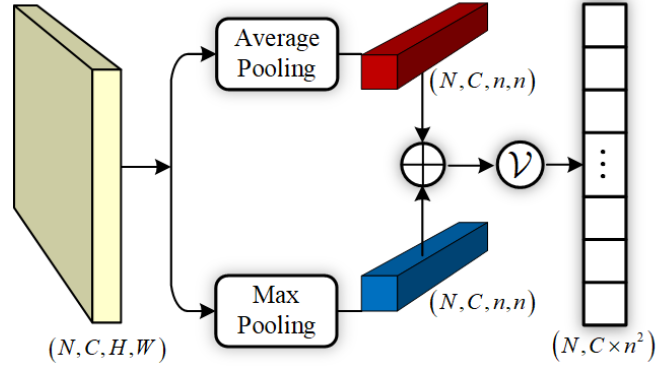$$\hat{f}_c = s_c \cdot f_c, \tag{2}$$



Fig. 3. The detailed process of the last feature map in each branch of HRNet-ReID. Both Adaptive Average Pooling(AAP) and Adaptive Max Pooling(AMP) are utilized to extract the channel features from the output of each branch. Then the two obtained features are added together, and further reshaped into a feature sequence.

where $f_c$ and $s_c$ denote the feature map and scaling factor in the same $c$-th channel of input image. The complete expression of $s$ is

$$s = S(W_U \delta(W_D z)), \tag{3}$$

where $W_U$ and $W_D$ denote the channel-upscaling layer and channel-downscaling layer with the ratio $r$ as the sampling rate, respectively. $S$ and $\delta$ denote Sigmoid gating [44] and the activation function ReLU [45], respectively. The channel-wise statistic $z \in \mathbb{R}^C$ is obtained from the feature map and the $c$-th element of $z$ is represented by the following formula:

$$z_c = H_{GP}(f_c) = \frac{1}{H \times W} \sum_{i=1}^{H} \sum_{j=1}^{W} f_c(i, j), \tag{4}$$

where $H_{GP}$ denotes global pooling function, and $f_c(i, j)$ denotes the pixel value at position $(i, j)$ of $c$-th feature $f_c$.

On the basis of RCAN's contribution, we adopt the VDSR as the super-resolution (SR) module which is fused with the channel attention (CA) mechanism to further enhance the performance. We add CA layer between every two layers of the original VDSR network. The fused SR module is named as VDSR-CA. The simplest Manhattan distance is utilized to train the VDSR-CA module, and the loss $\mathcal{L}_{SR}$ is calculated as:

$$\mathcal{L}_{SR} = \sum_{i=1}^{P \times K} ||\mathbb{E}_{x_l^i \sim D_l, x_h^i \sim D_h}[\mathcal{G}(x_l^i) - x_h^i]||_1, \tag{5}$$

where $\mathcal{G}$ denotes VDSR-CA module. $P$ and $K$ denote the number of selected persons and the number of corresponding images of each selected person, respectively.

### C. Innovation of High-Resolution Network

HRNet aims to fully extract features of input person images for retrieving and matching. Existing representation head of HRNet, such as HRNet-W32-C cannot perform satisfactorily in re-ID task. For this reason, we design a new representation head of HRNet which is adapted to re-ID task. We name the improved HRNet as HRNet-ReID.

Fig. 3 illustrates the elaborate processing of the last feature map in each branch. For the reason that feature maps with higher resolutions may contain more pixel space information,
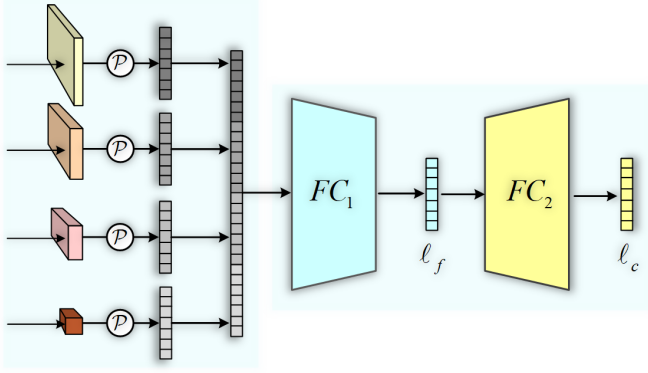
Fig. 4. Illustration of the proposed representation head in HRNet-ReID. For each branch, the output is first processed by the mapping $\mathcal{P}$ which is detailed introduced in Fig. 3. The extracted feature sequences $Seq^{(1\sim4)}$ represented by four short columns of different gray levels are concatenated together and denoted as the longest $Seq^{(5)}$. Then we obtain the feature representation $\ell_f$ and the classification output $\ell_c$ through two full connection layers $FC_1$ and $FC_2$.

we design a multi-resolution feature fusion strategy in HRNet. Adaptive Average Pooling (AAP) and Adaptive Max Pooling (AMP) operations are used to compress and refine the feature map information. For each branch, the two feature maps extracted by pooling layers are added into one feature map. Then the final sequence of each branch is obtained by reshaping the generated feature map. Specifically, the output sequence $Seq^{(n)}$ is defined as:

$$
\begin{aligned}
Seq^{(n)} &= \mathcal{P}(f^{(n)}, n; AAP, AMP, \mathcal{V}) \\
&= \mathcal{V}(AAP(f^{(n)}, n, n) + \lambda_n AMP(f^{(n)}, n, n)),
\end{aligned}
\tag{6}
$$

where $\mathcal{P}$ denotes the mapping from feature map to feature sequence. $f^{(n)}$ denotes the final feature map of HRNet in each branch. $n \in \{1, 2, 3, 4\}$ denotes the index of branch in HRNet, and $\mathcal{V}$ denotes the reshape operation which can transfer a feature map into a feature sequence. In particular, $f^{(n)}$ also represents the feature map containing the $n-$th highest resolution in the four branches, and the parameter $n$ also controls the output size of adaptive pooling layers.

Guided by the above designment, we use map $\mathcal{P}$ to serialize the four feature maps of HRNet and reconstruct the feature representations.

Fig. 4 shows the architecture of the entire output representation which contains $Seq^{(1\sim4)}$ obtained from map $\mathcal{P}$ and $Seq^{(5)}$ concatenated by $Seq^{(1\sim4)}$. We exploit a simple classification module to further process $Seq^{(5)}$ which consists of two fully connected layers denoted as $FC_1$ and $FC_2$. The outputs of them are represented as $\ell_f$ and $\ell_c$ respectively.

### D. HRNet-ReID Training

In recent years, a unified learning strategy which combines metric learning and representation learning has been widely applied in the person re-ID task and achieved great performance. In the training phase, the extracted $Seq^{(1\sim5)}$ in HRNet-ReID participate in metric learning, and the classification output $\ell_c$ participates in the representation learning. In the

testing phase, we concatenate $Seq^{(5)}$ and $\ell_f$ as the final feature representation for evaluation.

A representative loss function in metric learning called triplet loss is usually used for fine-grained recognition at the individual level. Here, we take advantage of the batch hard triplet loss $\mathcal{L}_{BH}$, which is calculated as:

$$
\begin{aligned}
\mathcal{L}_{BH} = \sum_{i=1}^{P} \sum_{a=1}^{K} \sum_{t=1}^{5} [m + \max_{p=1\ldots K}||Seq_{a,i}^{(t)} - Seq_{p,i}^{(t)}||_2 \\
- \min_{\substack{j=1\ldots P \\ n=1\ldots K \\ j \neq a}} ||Seq_{a,i}^{(t)} - Seq_{n,j}^{(t)}||_2]_+,
\end{aligned}
\tag{7}
$$

where $a$ denotes an anchor image, $p$ denotes a positive sample image, $n$ denotes a negative sample image, and $m$ represents a margin parameter to control the differences between intra and inter distances.

Moreover, we adopt the cross entropy label smooth loss that combines a label smoothing mechanism as the classification loss $\mathcal{L}_{CE}$ for representation learning:

$$
\mathcal{L}_{CE} = \sum_{n=1}^{P \times K} [-\sum_{y=1}^{M} \log(p(y))q(y)],
\tag{8}
$$

where $M$ represents the number of person labels involved in training set, and $p(y)$ denotes the probability that predicted label is $y$. Besides, the definition of $q(y)$ is:

$$
q(y) = \begin{cases} 1 - \frac{M-1}{M}\varepsilon & if \quad y = y_{truth} \\ \frac{\varepsilon}{M} & others \end{cases}
\tag{9}
$$

where $y_{truth}$ is the ground-truth label of the input image and $\varepsilon$ is a parameter of $\mathcal{L}_{CE}$.

Based on the two loss functions mentioned above (Eq. (7) and (8)), for each batch of training set, we compute the HRNet-ReID loss $\mathcal{L}_{ID}$ by:

$$
\mathcal{L}_{ID} = \lambda_{CE}\mathcal{L}_{CE} + \lambda_{BH}\mathcal{L}_{BH},
\tag{10}
$$

where $\lambda_{CE}$ and $\lambda_{BH}$ are parameters to control the importance of $\mathcal{L}_{CE}$ and $\mathcal{L}_{BH}$, respectively.

It is worth emphasizing that, in the whole architecture of our PS-HRNet, HRNet-ReID as a feature extraction network, which plays different roles in different phases. As shown in Fig. 2, HRNet-ReID undertakes the task of effectively extracting discriminating features from input HR images in the first phase of the blue arrow representation. In the second phase of the black arrow representation, HRNet-ReID of second phase both participates in multi-task joint learning and forms a pseudo-siamese framework with the first phase HRNet-ReID. For the convenience of representation and distinction, when the resolution of the input images is HR or LR, we define HRNet-ReID as HRNet-ReID-H and HRNet-ReID-L, respectively.

### E. Multi-Task Learning Under Pseudo-Siamese Framework

The design and application of the pseudo-siamese framework make the training mode of PS-HRNet different from most other cross-resolution person re-ID methods. To build a joint multi-task learning structure as the black arrow shown in Fig. 2, we specially design a set of training strategies.

In the first phase, single HRNet-ReID is trained individually with the HR images from $\mathcal{D}_h$. By minimizing the loss $\mathcal{L}_{ID}$ and observing the indicators of HRNet-ReID in the testing set,

---

**Algorithm 1 PS-HRNet module training**

---

**Input:** Training set $\mathcal{D}_h = \{x_h^i, y^i\}_{i=1}^N$ and $\mathcal{D}_l = \{x_l^i, y^i\}_{i=1}^N$.

**Output:** Joint re-ID module composed of **VDSR-CA** and **HRNet-ReID-L**.

**Phase 1 (Preparation)**
    Take $\mathcal{D}_h$ as input.
    **for** $i = 1$ **to** $iter_1$ **do**
        Update the single HRNet-ReID-H with the loss $\mathcal{L}_{ID}$ (Eq. (10)).
    **end for**

**Phase 2 (Joint Multi-Task Learning)**
    Take $\mathcal{D}_h$ and $\mathcal{D}_l$ as input.
    Import the first phase trained HRNet-reID-H module.
    **for** $i = 1$ **to** $iter_2$ **do**
        Update the joint multi-task learning loss with the loss $\mathcal{L}_{TOTAL}$
        (Eq. (12)).
    **end for**

---

we obtain a high-performance HRNet-ReID module defined as HRNet-ReID-H for subsequent joint learning and construction of the pseudo-siamese framework.

In the second phase, we first concatenate VDSR-CA with the newly defined HRNet-ReID-L, and then construct the pseudo-siamese framework with the obtained HRNet-ReID-H from the first phase and HRNet-ReID-L from the concatenated structure via the Manhattan distance loss $\mathcal{L}_{\mathcal{PS}}$. The pseudo-siamese framework is adopted to measure the similarity of two inputs. Compared with the siamese framework, the pseudo-siamese framework is more suitable for the situation where two inputs have a certain difference. We first define the ordered set $A = \left\{ Seq^{(1)}, Seq^{(2)}, Seq^{(3)}, Seq^{(4)}, Seq^{(5)}, \ell_c \right\}$. Then the loss $\mathcal{L}_{\mathcal{PS}}$ is defined as:

$$\mathcal{L}_{\mathcal{PS}} = \sum_{i=1}^m ||C_h^{(i)} - C_l^{(i)}||_1, \qquad (11)$$

where ordered set $C \subseteq A$ and $C \neq \varnothing$, denotes a combination of elements in ordered set $A$. $C_h$ and $C_l$ denote the set of elements from HRNet-ReID-H and HRNet-ReID-L under the same combination, respectively. $m$ denotes the amount of elements in ordered set $C$. The process of joint multi-task learning depends on the training set from both $\mathcal{D}_h$ and $\mathcal{D}_l$. The VDSR-CA module first reads the LR images from $\mathcal{D}_l$ for super-resolution reconstruction, and outputs the restored images with the same resolution as the HR images. The restored images are read by HRNet-ReID-H to train itself with the loss $\mathcal{L}_{ID}$ (Eq. (10)), and at the same time, they participate in the calculation of $\mathcal{L}_{SR}$ (Eq. (5)) together with the corresponding HR images from $\mathcal{D}_h$. Finally, the guidance of HRNet-ReID-H to HRNet-ReID-L is realized by loss $\mathcal{L}_{\mathcal{PS}}$ (Eq. (11)).

To sum up, we combine the ID loss (Eq. (10)), the SR loss (Eq. (5)) and the PS loss (Eq. (11)) into the joint multi-task learning loss $\mathcal{L}_{TOTAL}$, defined as:

$$\mathcal{L}_{TOTAL} = \mathcal{L}_{ID} + \lambda_{SR}\mathcal{L}_{SR} + \lambda_{PS}\mathcal{L}_{PS}, \qquad (12)$$

where $\lambda_{SR}$ and $\lambda_{PS}$ are two parameters to control the importance of $\mathcal{L}_{SR}$ and $\mathcal{L}_{PS}$, respectively. Algorithm 1 gives a summary of the entire training process.

When in the testing phase, we uniformly input the gallery and query images into the trained network composed of VDSR-CA and HRNet-ReID-L. The generated $Seq^{(5)}$ and $\ell_f$ are concatenated as the ultimate feature representation for matching and evaluation. The entire testing process is executed end-to-end without additional operations.

## IV. EXPERIMENTS

In this part, we first give detailed descriptions of the datasets for evaluation, the experimental settings, and the specific implementation details. Then, a number of experiments are carried out to prove the validity of our proposed method through the comparison with existing methods and ablation studies.

### A. Datasets

Our experiment involves nine datasets, including four high-resolution re-ID datasets for traditional re-ID task: VIPeR [46], CUHK03 [47], DukeMTMC-reID [48], Market-1501 [49], and four synthetic cross-resolution datasets named as MLR datasets which are constructed from traditional versions: MLR-VIPeR, MLR-CUHK03, MLR-DukeMTMC-reID, MLR-Market1501, as well as one native cross-resolution dataset sampled in real world: CAVIAR [50]. The five cross-resolution datasets involved in our experiments are described as follows:

*1) MRL-VIPeR:* The MLR-VIPeR includes 632 person-image pairs captured from 2 cameras, a total of 1264 pictures. Following [16], we randomly down-sample all the pictures captured by one of the cameras by a down-sampling rate $r \in \{2, 3, 4\}$ and the images collected from another camera remain unchanged. Here we use the standard 316/316 training/testing identity split.

*2) MLR-CUHK03:* The MLR-CUHK03 dataset is generated from images taken by 10 (5 pairs) different cameras. It contains 14097 images of 1467 individuals. As [16], for each pair of cameras, we down-sample the images captured from one camera by randomly selecting a down-sampling rate $r \in \{2, 3, 4\}$, while the resolution of the images collected by the other cameras are unchanged. Here we use the 1,367/100 training/testing identity split.

*3) MLR-DukeMTMC-reID:* The MLR-DukeMTMC dataset includes 36,411 images of 1,404 identities captured from 8 cameras. Following [39], we randomly select one camera, and down-sample the images by the same down-sampling rate while the image resolution of other cameras remains unchanged. We use the standard 702/702 training/testing identity split.

*4) MLR-Market-1501:* The MLR-Market-1501 dataset is composed of 32,668 pictures of 1,501 persons captured by 6 cameras. Following [39], we preprocess images of one camera with the same down-sampling rate, while other image resolutions remain unchanged. According to the person ID label, the dataset is separated into a training set containing 751 pedestrians and a testing set containing 750 pedestrians.

*5) CAVIAR:* The CAVIAR is a dataset collected in real world. It consists of 1,220 pictures of 72 persons collected from 2 cameras. Following [16], 22 persons are discard by us with only HR images and we randomly split the dataset into two halves based on 25 identities labels for training and testing, respectively.

TABLE I
EXPERIMENTAL RESULTS OF CROSS-RESOLUTION PERSON RE-ID (%). THE BOLD AND UNDERLINED NUMBERS INDICATE TOP TWO RESULTS, RESPECTIVELY.

| Module | MLR-Market-1501 | | | MLR-CUHK03 | | | MLR-VIPeR | | | MLR-DukeMTMC-reID | | | CAVIAR | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Rank1 | Rank5 | Rank10 | Rank1 | Rank5 | Rank10 | Rank1 | Rank5 | Rank10 | Rank1 | Rank5 | Rank10 | Rank1 | Rank5 | Rank10 |
| PCB [51] | 76.9 | 88.9 | 92.4 | 75.3 | 92.7 | 98.1 | 42.6 | 65.8 | 75.9 | 66.4 | 82.5 | 87.1 | 35.9 | 72.1 | 88.6 |
| DenseNet-121 [52] | 60 | 78.8 | - | 70.8 | 91.3 | - | 31.4 | 63.1 | - | - | - | - | 31.1 | 65.5 | - |
| ResNet-50 [53] | 57 | 78.7 | - | 67.4 | 91.7 | - | 29.9 | 62.2 | - | - | - | - | 29.6 | 64 | - |
| SE-ResNet-50 [54] | 58.2 | 78.6 | - | 70.8 | 92.3 | - | 33.5 | 63.6 | - | - | - | - | 30.8 | 65.1 | - |
| SPreID [28] | 77.4 | 89 | 93.9 | 76.5 | 92.5 | 98.3 | 42.4 | 65.8 | 75.1 | 68.4 | 84.5 | 89.1 | 36.2 | 71.9 | 88.7 |
| Part Aligned [55] | 75.6 | 88.5 | 92.2 | 73.4 | 92.1 | 97.5 | 40.2 | 62.3 | 73.1 | 67.5 | 83.1 | 87.2 | 35.7 | 71.4 | 87.9 |
| CamStyle [56] | 74.5 | 88.6 | 93 | 69.1 | 89.6 | 93.9 | 34.4 | 56.8 | 66.6 | 64 | 78.1 | 84.4 | 32.1 | 72.3 | 85.9 |
| PyrNet [57] | 83.8 | 93.3 | 95.6 | 83.9 | 97.1 | 98.5 | - | - | - | 79.6 | 88.1 | 91.2 | 43.6 | 79.2 | 90.4 |
| FD-GAN [58] | 79.6 | 91.6 | 93.5 | 73.4 | 93.8 | 97.9 | 39.1 | 62.1 | 72.5 | 67.5 | 82 | 85.3 | 33.5 | 71.4 | 86.5 |
| JUDEA [14] | - | - | - | 26.2 | 58 | 73.4 | 26 | 55.1 | 69.2 | - | - | - | 22 | 60.1 | 80.8 |
| SDF [59] | - | - | - | 22.2 | 48 | 64 | 9.3 | 38.1 | 52.4 | - | - | - | 14.3 | 37.5 | 62.5 |
| SLD$^2$L [13] | - | - | - | - | - | - | 20.3 | 44 | 62 | - | - | - | 18.4 | 44.8 | 61.2 |
| SING [16] | 74.4 | 87.8 | 91.6 | 67.7 | 90.7 | 94.7 | 33.5 | 57 | 66.5 | 65.2 | 80.1 | 84.8 | 33.5 | 72.7 | 89 |
| CSR-GAN [38] | 76.4 | 88.5 | 91.9 | 71.3 | 92.1 | 97.4 | 37.2 | 62.3 | 71.6 | 67.6 | 81.4 | 85.1 | 34.7 | 72.5 | 87.4 |
| FFSR [60] | 59.2 | 80.1 | - | 70.5 | 92.3 | - | 40.3 | 65.3 | - | - | - | - | 31.1 | 68.7 | - |
| RIFE [60] | 62.6 | 82.4 | - | 69.7 | 91.5 | - | 33.9 | 63.6 | - | | | | 35.7 | 74.9 | - |
| FFSR+RIFE [60] | 66.9 | 84.7 | | 73.3 | 92.6 | - | 41.6 | 64.9 | - | | | | 36.4 | 72 | |
| RAIN [15] | - | - | - | 78.9 | 97.3 | 98.7 | 42.5 | 68.3 | 79.6 | - | - | - | 42 | 77.3 | 89.6 |
| CAD-Net [39] | 83.7 | 92.7 | 95.8 | 82.1 | 97.4 | 98.8 | 43.1 | 68.2 | 77.5 | 75.6 | 86.7 | 89.6 | 42.8 | 76.2 | 91.5 |
| CAD-Net++ [61] | 84.1 | 93 | 96.2 | 83.4 | 98.1 | 99.1 | 43.4 | 68.7 | 78.2 | 77.2 | 88.1 | 90.4 | 43.1 | 76.5 | 92.3 |
| INTACT [40] | <u>88.1</u> | <u>95</u> | 96.9 | 86.4 | 97.4 | 98.5 | <u>46.2</u> | <u>73.1</u> | <u>81.6</u> | 81.2 | 90.1 | 92.8 | 44 | 81.8 | 93.9 |
| PRI [62] | 84.9 | 93.5 | 96.1 | 85.2 | 97.5 | 98.8 | - | - | - | 78.3 | 87.5 | 91.4 | 43.2 | 78.5 | 91.9 |
| PCB+PRI [62] | <u>88.1</u> | 94.2 | 96.5 | 86.2 | <u>97.9</u> | <u>99.1</u> | - | - | - | 81.6 | 89.6 | <u>92.4</u> | 44.3 | 83.7 | <u>94.8</u> |
| PyrNet+PRI [62] | 86.9 | 93.8 | 96.4 | <u>86.5</u> | 97.7 | <u>99.1</u> | - | - | - | <u>82.1</u> | **91.1** | 92.8 | <u>45.2</u> | <u>84.1</u> | 94.6 |
| **Ours** | **91.5** | **96.7** | **97.9** | **92.6** | **98.3** | **99.4** | **48.7** | **73.4** | **81.7** | **82.3** | <u>90.5</u> | 92.8 | **48.2** | **84.5** | **96.3** |

## B. Experimental Settings

In the training phase, both the four traditional re-ID datasets and the corresponding cross-resolution datasets are applied to train the module. The training process is divided into two phases: In the first phase, we only train the HRNet-ReID-H on the traditional high-resolution re-ID datasets, and obtain a high-performance module; In the second phase, we combine the VDSR-CA module with the HRNet-ReID-L guided by HRNet-ReID-H and jointly train them on both traditional datasets and cross-resolution datasets.

In the testing phase, we evaluate the performance of our proposed PS-HRNet on five cross-resolution datasets, where the query sets contain LR images and the gallery sets contain HR images. In particular, since CAVIAR is a genuine cross-resolution dataset, we follow the experimental setting in [16] for training and testing. We adopt the standard single-shot person re-ID settings, and apply the average cumulative match characteristic (CMC) to quantify the performance and report the results of ranks 1, 5 and 10.

## C. Implementation details

In the VDSR-CA module, we retain the entire structure of the original VDSR network, and embed the Channel Attention

(CA) mechanism proposed in the RCAN [43]. The internal parameter $r$ of CA mechanism is set to 4.

With respect to HRNet-ReID module, we select the HRNet-W32-C pretrained with the ImageNet dataset as the backbone for feature extraction. Here we redesign the classifier to adapt to the re-ID task as shown in Fig. 4. The lengths of $Seq^{(1\sim5)}$ are set to 2048, 1024, 2048, 1024 and 6144, respectively. Besides, the dimension of $\ell_f$ is set to 512, and the dimension of $\ell_c$ is equal to the number of target categories.

Before training, all images are resized to $256 \times 128 \times 3$. A mini-batch contains 24 pairs of images of $P = 4$ persons, and each person has $K = 6$ pairs of HR and LR images. We choose SGD to optimize our module with weight decay $5 \times 10^{-4}$. The learning rates for training HRNet-ReID and VDSR-CA are set to $8.5 \times 10^{-3}$ and $8.5 \times 10^{-4}$, respectively, which are decreased by 0.1 every 30 epochs. Our module is trained for 70 epochs in total. The hyper-parameters $m$, $\lambda_{CE}$, $\lambda_{BH}$, $\lambda_{SR}$ and $\lambda_{PS}$ are set to 0.1, 1.15, 0.2, 0.5 and 0.5, respectively. In $\mathcal{L}_{PS}$, we select $\{Seq^{(1)}, Seq^{(4)}, Seq^{(5)}, \ell_c\}$ as a combination to participate in the operation. Some data augmentation tricks are utilized, such as random flipping, padding and random cropping. We perform our experiments with PyTorch of version 1.6 on single 11GB NVIDIA RTX 2080Ti GPU.

TABLE II
EVALUATING OURS LOSS COMPONENTS ON MLR-MARKET1501. ID: IDENTITY CLASSIFICATION LOSS (EQ. (10)), SR: SUPER-RESOLUTION LOSS (EQ. (5)), PS: OUR PSEUDO-SIAMESE FRAMEWORK LOSS (EQ. (11)).

| Supervision | Rank1 | Rank5 | Rank10 |
|---|---|---|---|
| ID | 84.7 | 92 | 94.6 |
| SR+ID | 87.6 | 95.1 | 97 |
| SR+ID+PS | 91.5 | 96.7 | 97.9 |

TABLE III
PERFORMANCE COMPARISON TEST OF SUPER-RESOLUTION RECONSTRUCTION AND CROSS-RESOLUTION PERSON REID ON THE MLR-CUHK03 TEST SET.

| Module | SSIM | PSNR | Rank1 |
|---|---|---|---|
| CycleGAN [63] | 0.55 | 14.1 | 62.1 |
| SING [16] | 0.65 | 18.1 | 67.7 |
| CSR-GAN [38] | 0.76 | 21.5 | 71.3 |
| CAD-Net [39] | 0.73 | 20.2 | 82.1 |
| Solo SR | 0.93 | 27.9 | 0.04 |
| SR+ID | 0.91 | 26.1 | 88.9 |
| SR+ID+PS | 0.91 | 26.1 | 92.6 |

## D. Comparisons to State-of-the-Art Methods

We compare our PS-HRNet with a series of far-ranging state-of-the-art person re-ID methods, which can be roughly separated into two main categories. (1) Conventional methods designed for traditional person re-ID task: PCB [51], DenseNet-121 [52], ResNet-50 [53], SE-ResNet-50 [54], SPreID [28], Part Aligned [55], CamStyle [56] and FD-GAN [58]; (2) Pointed methods designed for cross-resolution person re-ID task: JUDEA [14], SDF [59], SLD$^2$L [13], SING [16], CSR-GAN [38], FFSR [60], RIFE [60], FFSR+RIFE [60], RAIN [15], CAD-Net [39], CAD-Net++ [61], PRI [62], PCB+PRI [62] and PyrNet+PRI [62]. These methods in the comparison almost cover all the current methods in the cross-resolution re-ID field.

The comparison results of the above approaches on five datasets are listed in Table I. We can evidently observe that:

• Our PS-HRNet obtains state-of-the-art performance on all five datasets, and its Rank-1 outperforms the best competitor by 6.1% on the MLR-CUHK03 dataset.

• Compared with the conventional person re-ID methods, our method outperforms their best result by 11.9% on the MLR-Market-1501 dataset, which indicates that the information hided in LR images cannot be extracted and utilized effectively by conventional methods when processing cross-resolution person images. Besides, it also demonstrates that the super-resolution reconstruction module with channel attention mechanism plays a significant role in the cross-resolution re-ID task.

• Compared with the pointed methods designed for solution of cross-resolution person re-ID problem, PS-HRNet outperforms all existing methods, which reflects the importance of the tailor-made high-resolution feature extraction network and pseudo-siamese framework.

## E. Ablation Study

1) Analysis of Loss Functions: Our PS-HRNet jointly trains image SR module, feature extraction network and pseudo-siamese framework with several loss functions (cf. Eq. (12) etc.). We use the same research strategy as described in INTACT [40] to study the effectiveness of different losses in PS-HRNet on the MLR-Market-1501 dataset. Table II reports the ablation results, which reflect that:

• When only ID loss is included, compared with Table I, it can be obviously found that the performance of our method surpasses almost all existing methods except INTACT

[40] and PCB+PRI [62] on the Rank-1, and even reaches the same performance level of single PRI [62]. The results reflect the great feature extraction ability of HRNet on cross-resolution person images which should be owed to its unique high-resolution parallel structure. This also demonstrates the necessity and rationality of applying HRNet to process low-resolution images.

• With the addition of SR loss, the performance is further improved by 2.9%, which can prove that the image reconstruction function provided by VDSR-CA module has a positive effect on cross-resolution person re-ID.

• The addition of pseudo-siamese framework loss ultimately increases the Rank-1 by 3.9%, which verifies the positive effect of the pseudo-siamese framework. Furthermore, it proves the necessity of reducing the discrepancies of feature space between HR and LR images, which has been overlooked all along.

• Following INTACT [40], experiments on loss are conducted to explore the influence of different loss combinations on image restoration quality and whether it will affect the restored images after using the pseudo-siamese framework. We use PSNR and SSIM which are two quantitative indicators that reflect the quality of image restoration to evaluation. We compare our method with CycleGAN [63], SING [16], CSR-GAN [38] and CAD-Net [39]. According to the comparative results in Table III, the solo SR module achieves the best performance in image super-resolution reconstruction which also confirms the effectiveness of VDSR-CA. After using SR+ID, the PSNR and SSIM are slightly reduced, but the recognition accuracy is greatly improved. Furthermore, after adding the pseudo-siamese framework, the value of PSNR and SSIM have no changed, which indicates that the pseudo-siamese framework does not affect the quality of the restored images on visual level.

• Following INTACT [40], Fig. 5 shows some person images restored by our method. These images are selected from the testing set of MLR-CUHK03. Although our PS-HRNet is superior to the existing baseline method in recognition accuracy, there is still a gap between the recovery quality of low-resolution images and the ground truth on visual level. Such outcome does make perfect sense. As we mentioned earlier, the focus is not on the restoration of image quality but on the examination of deep-level feature and
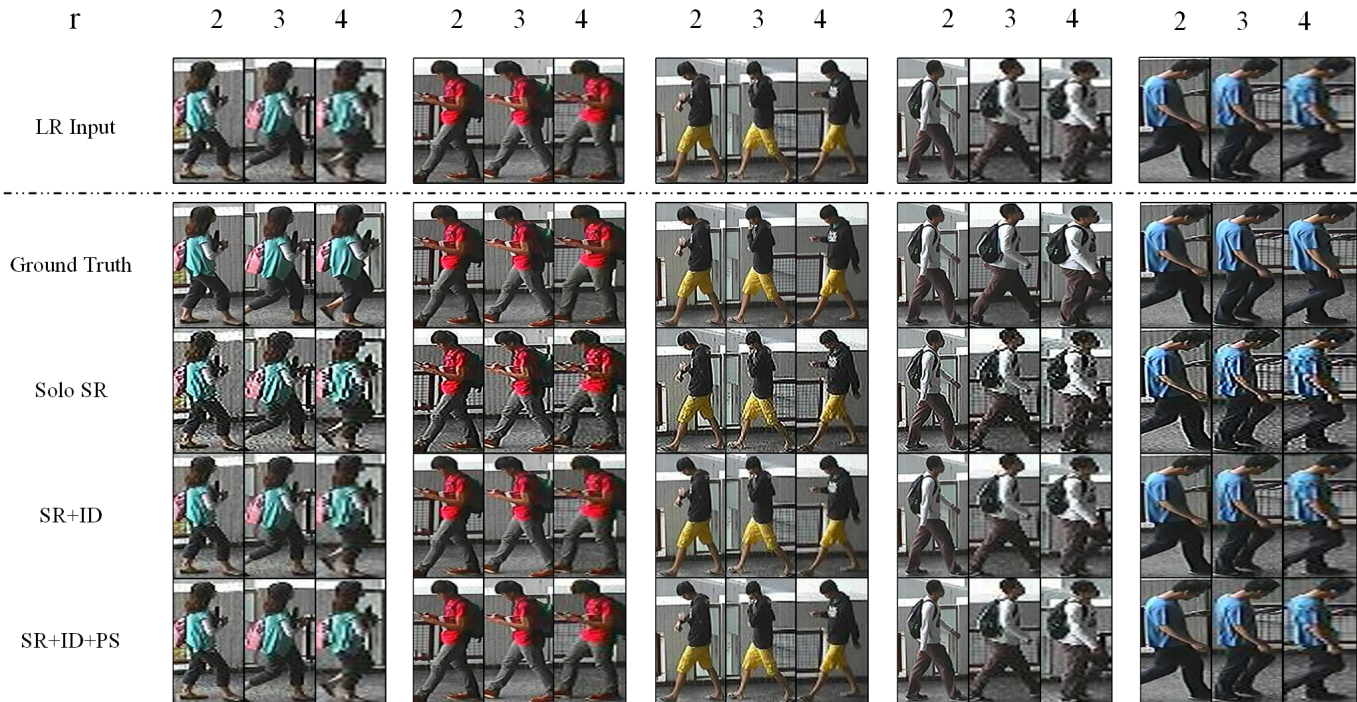
Fig. 5. Visual comparisons of restored HR test images from the MLR-CUHK03 dataset. Given input images of 3 different low-resolutions $r \in 2, 3, 4$, our PS-HRNet outputs the corresponding restored HR images of solo SR module, SR+ID module and SR+ID+PS module.

### TABLE IV
RECOGNITION ACCURACY (%) OF DIFFERENT SR MODULES ON MLR-MARKET1501.

| Module | Rank1 | Rank5 | Rank10 |
|---|---|---|---|
| Bicubic +Resnet50 | 64.7 | 82.1 | 89.1 |
| SRCNN [19] +Resnet50 | 62.1 | 80.9 | 87.6 |
| VDSR [22] +Resnet50 | 68.2 | 86.5 | 90.7 |
| RCAN [43] +Resnet50 | 73.1 | 89.6 | 92.4 |
| VDSR-CA+Resnet50 | 72.3 | 88.2 | 92 |

### TABLE V
RECOGNITION ACCURACY (%) OF DIFFERENT POOLING METHOD ON MARKET1501 AND DUKEMTMC-REID.

| Method | Market-1501 | | | | DukeMTMC | | | |
|---|---|---|---|---|---|---|---|---|
| | Rank1 | Rank5 | Rank10 | mAP | Rank1 | Rank5 | Rank10 | mAP |
| AAP | 94 | 97.4 | 98.1 | 83.8 | 86.4 | 93.2 | 95.2 | 73.6 |
| AMP | 93.6 | 97.7 | 98.7 | 83.4 | 86.6 | 93.7 | 95.4 | 74 |
| AAP+AMP | 94.4 | 98.2 | 98.6 | 84.1 | 87 | 94.5 | 96 | 74.3 |

### TABLE VI
RECOGNITION ACCURACY (%) OF DIFFERENT FEATURE EXTRACTION NETWORKS ON MARKET1501 AND DUKEMTMC-REID.

| Module | Backbone | Market-1501 | | DukeMTMC-reID | |
|---|---|---|---|---|---|
| | | Rank1 | mAP | Rank1 | mAP |
| DPFL [21] | ResNet-50 | 88.9 | 73.1 | 79.2 | 60.6 |
| PCB [51] | ResNet-50 | 92.3 | 77.4 | 81.7 | 66.1 |
| FD-GAN [58] | ResNet-50 | 90.5 | 77.7 | 80 | 64.5 |
| CamStyle [56] | ResNet-50 | 89.2 | 71.6 | 78.6 | 57.6 |
| DenseNet-121 [52] | DenseNet | 90.2 | 74 | 67.4 | 46.2 |
| HRNet-W32-C [64] | HRNetV2-W32 | 87.8 | 72.8 | 78.6 | 60.3 |
| HRNet-ReID | HRNetV2-W32 | 94.4 | 84.1 | 87 | 74.3 |

semantic information. Blindly using visual sensory experience and related indicators as a standard for image restoration is of little significance because the way the neural network observes the image is completely different from the human visual mechanism. This also gives a reasonable interpretation on the variation of PSNR and SSIM in Table III. Therefore, we do not recommend subsequent studies on cross-resolution person re-identification to conduct excessive super-resolution reconstruction experiments at the aspect of human vision.

*2) Analysis of VDSR-CA:* In order to further explore the effectiveness of the VDSR-CA, we adopt ResNet-50 [53] as a unified feature extraction network on the MLR-Market1501 dataset, and test the effect of using cubic interpolation, SR-CNN [19], VDSR-CA and RCAN [43] as super-resolution modules for joint learning.

The recognition accuracy of different SR modules on MLR-Market1501 are presented in Table IV, which reflects that the performance of our VDSR-CA module goes far beyond the classic SRCNN network and is better than the original VDSR

network. Although the Rank indicators are slightly lower than the result of joint learning with RCAN, VDSR-CA has a lighter architecture than RCAN which means more practical in real-world applications. Considering the above factors, our VDSR-CA is the best choice for super-resolution modules.

TABLE VII
RECOGNITION ACCURACY (%) OF DIFFERENT SEQUENCE COMBINATIONS
ON MLR-MARKET1501.

| Sequences | Rank1 | Rank5 | Rank10 |
|---|---|---|---|
| $\{Seq^{(5)}, \ell_c\}$ | 90.5 | 96.6 | 97.9 |
| $\{Seq^{(1)}, Seq^{(2)}, Seq^{(3)}, Seq^{(4)}, Seq^{(5)}, \ell_c\}$ | 90.9 | 96.1 | 97.5 |
| $\{Seq^{(1)}, Seq^{(4)}, Seq^{(5)}, \ell_c\}$ | 91.5 | 96.7 | 97.9 |

*3) Analysis of HRNet-ReID:* The function of HRNet-ReID aims to fully extract features of input person images for retrieving and matching. In order to make a large number of high-dimensional feature map data extracted by the previous HRNet-W32 can be effectively processed by the classification layer at the end of HRNet-ReID, we introduce adaptive average pooling and adaptive max pooling as important means of feature compression.

Table V demonstrates the effect of different pooling options on the recognition accuracy. Obviously, under the same experimental conditions, only the simultaneous application of AAP and AMP can obtain higher detection accuracy. Theoretically, maximum pooling and average pooling play specific roles in extracting feature texture information and global background information respectively. Therefore, the combination of average pooling and adaptive max pooling is valid.

For further investigating the performance of HRNet-ReID as the backbone in the person re-ID task, we conduct experiments on the following conventional person re-ID network combined with ResNet-50, PCB and DenseNet as the comparisons with HRNet-ReID on the Market-1501 and DukeMTMC-reID datasets: DPFL [21], PCB [51], FD-GAN [58], CamStyle [56], DenseNet121 [52]. In addition, in order to validate the effectiveness of our designed HRNet representation head, we use the original version of HRNet-W32-C [64] as reference. The representation head of HRNet-W32-C is proposed to solve the problem of image classification on the ImageNet dataset and achieves good performance.

As shown in Table VI, the experimental results reflect that our HRNet-ReID significantly outperforms the original HRNet-W32-C and other methods on both datasets which confirms the effectiveness of our modified representation head.

*4) Analysis of pseudo-siamese framework:* For the purpose of exploring the impact of different training strategies under the pseudo-siamese framework, we select multiple sets of $Seq^{(n)}$ and $\ell_c$ to participate in the training of pseudo-siamese framework loss and perform testing on the MLR-Market-1501 dataset. According to the combinatorial mathematics, there are 84 combinations for the calculation of Eq. (11) in theory. Here we select three typical combinations for evaluation.

The experimental results are listed in Table VII. We can clearly observe that different selections of sequences bring different performances. Here we adopt the combination of $\{Seq^{(1)}, Seq^{(4)}, Seq^{(5)}, \ell_c\}$ in our training strategy. Limited by time, there may be a better combination that further improves the performance of our PS-HRNet.

## V. CONCLUSION

In this article, we design a novel approach named Deep High-Resolution Pseudo-Siamese Framework (PS-HRNet) to significantly alleviate the the resolution mismatch problem and improve recognition accuracy in cross-resolution person re-ID task. Our framework utilizes VDSR-CA as the super-resolution module, HRNet-ReID as the feature extraction network. The former integrates channel attention mechanism into VDSR, which can reasonably utilize the valuable high frequency components contained in different channels of feature map and restore the missing discriminative information in LR images effectively. The latter possesses a novel representation head designed by us which can effectively extract fine-grained details from cross-resolution person images. What's more, the pseudo-siamese framework is adopted and plays a significant component in reducing the distribution difference in feature information between LR and HR images. With extensive experiments, the results confirm that our PS-HRNet can extract discriminating and robust feature representations from cross-resolution images, and achieves the state-of-the-art performance in existing five benchmarks.

## REFERENCES

[1] C. Liu, S. Gong, C. C. Loy, and X. Lin, "Person re-identification: What features are important?" *in Proc. of the European Conf. Comput. Vis.* Springer, 2012, pp. 391–401.

[2] R. Layne, T. M. Hospedales, S. Gong, and Q. Mary, "Person re-identification by attributes," *in Proc. of the British Machine Vis. Conf.*, vol. 2, no. 3, 2012, p. 8.

[3] G. Wang, Y. Yuan, J. Li, S. Ge, and X. Zhou, "Receptive Multi-Granularity Representation for Person Re-Identification," *IEEE Trans. Image Processing*, vol. 29, pp. 6096-6109, 2020.

[4] K. Wang, C. Ding, S. Maybank, and D. Tao, "CDPM: Convolutional Deformable Part Models for Semantically Aligned Person Re-Identification," *IEEE Trans. Image Processing*, vol. 29, pp. 3416–3428, 2020.

[5] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," *in Proc. of the IEEE Int. Conf.Comput. Vis*, 2015, pp. 1116–1124.

[6] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification:Past, present and future," *arXiv preprint arXiv:1610.02984*, 2016.

[7] Z. Zhang, Y. Xie, W. Zhang, Y. Tang, and Q. Tian, "Tensor Multi-Task Learning for Person Re-Identification," *IEEE Trans. Image Processing*, vol. 29, pp. 2463-2477, 2020.

[8] Y. Zhai, S. Lu, Q. Ye, X. Shan, J. Chen, Y. Tian, and R. Ji, "Ad-cluster: Augmented discriminative clustering for domain adaptive person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, pp. 9021–9030, 2020.

[9] D. Chung, K. Tahboub, and E. J. Delp, "A two stream siamese convolutional neural network for person re-identification," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, 2017, pp. 1983–1991.

[10] Z. Zhang, T. Si, and S. Liu, "Integration convolutional neural network for person re-identification in camera networks," *IEEE Access*, vol. 6,pp. 36 887–36 896, 2018.

[11] Y. Wu, O. Bourahla, X. Li, F. Wu, Q. Tian, and X. Zhou, "Adaptive Graph Representation Learning for Video Person Re-Identification," *IEEE Trans. Image Processing*, vol. 29, pp. 8821–8830, 2020.

[12] Q. Zhou, B. Zhong, X. Lan, G. Sun, Y. Zhang, B. Zhang and R. Ji, "Fine-Grained Spatial Alignment Model for Person Re-Identification With Focal Triplet Loss," *IEEE Trans. Image Processing*, vol. 29, pp. 7578-7589, 2020.

[13] X.-Y. Jing, X. Zhu, F. Wu, X. You, Q. Liu, D. Yue, R. Hu, and B. Xu, "Super-resolution person re-identification with semi-coupled low-rank discriminant dictionary learning," *in Proc. of the IEEE Conf. Comput.Vis. and Pattern Recognition (CVPR)*, June 2015.

[14] X. Li, W.-S. Zheng, X. Wang, T. Xiang, and S. Gong, "Multi-scale learning for low-resolution person re-identification," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, December 2015.

[15] Y.-C. Chen, Y.-J. Li, X. Du, and Y.-C. F. Wang, "Learning resolution-invariant deep representations for person re-identification," *in Proc. of the AAAI Conf. Artificial Intelligence*, vol. 33, 2019, pp. 8215–8222.

[16] J. Jiao, W.-S. Zheng, A. Wu, X. Zhu, and S. Gong, "Deep low-resolution person re-identification," *in Proc. of the AAAI Conf. Artificial Intelligence*, vol. 32, no. 1, 2018.

[17] G. Z. A, J. Y. A, Y. Z. A, Z. L. A, and J. Z. B, "Optimal discriminative feature and dictionary learning for image set classification," *Information Sciences*, vol. 547, pp. 498–513, 2021.

[18] G. Zhang, H. Sun, F. Porikli, Y. Liu, and Q. Sun, "Optimal couple projections for domain adaptive sparse representation-based classification," *IEEE Trans. Image Processing*, vol. PP, no. 12, pp. 1–1, 2017.

[19] C. Dong, C. C. Loy, K. He, and X. Tang, "Learning a deep convolutional network for image super-resolution," *in Proc. of the European Conf. Comput. Vis.* Springer, 2014, pp. 184–199.

[20] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang et al., "Photo-realistic single image super-resolution using a generative adversarial network," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2017, pp. 4681–4690.

[21] 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 Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2016, pp. 1335–1344.

[22] J. Kim, J. Kwon Lee, and K. Mu Lee, "Accurate image super-resolution using very deep convolutional networks," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2016, pp. 1646–1654.

[23] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep high-resolution representation learning for human pose estimation," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2019, pp. 5693–5703.

[24] J. Bromley, J. W. Bentz, L. Bottou, I. Guyon, Y. LeCun, C. Moore, E. Säckinger, and R. Shah, "Signature verification using a "siamese" time delay neural network," *Int. Journal of Pattern Recognition and Artificial Intelligence*, vol. 7, no. 04, pp. 669–688, 1993.

[25] M. Zheng, S. Karanam, and Z. Wu, R. Radke, "Re-Identification with Consistent Attentive Siamese Networks," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2019, pp. 5735-5744.

[26] J. Liu, B. Ni, Y. Yan, P. Zhou, S. Cheng, and J. Hu, "Pose transferrable person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2018, pp. 4099–4108.

[27] H. Li, K. Qiu, L. Chen, X. Mei, L. Hong, and C. Tao, "Scattnet: Semantic segmentation network with spatial and channel attention mechanism for high-resolution remote sensing images," *IEEE Geoscience and Remote Sensing Letters*, 2020.

[28] M. M. Kalayeh, E. Basaran, M. Gökmen, M. E. Kamasak, and M. Shah, "Human semantic parsing for person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2018, pp. 1062–1071.

[29] M. Jia, X. Cheng, Y. Zhai, S. Ma, J. Zhang, and S. Lu, "Matching on Sets: Conquer Occluded Person Re-Identification Without Alignment," *in Proc. of the AAAI Conf. Artificial Intelligence*, 2021.

[30] G. Wang, S. Yang, H. Liu, Z. Wang, Y. Yang, S. Wang, Y. Gang, E. Zhou, and J. Sun, "High-order information matters: Learning relation and topology for occluded person re-identification," *in Proc. of IEEE Conf. Comput. Vis. and Pattern Recognition*, 2020, pp. 6449–6458.

[31] H. Wang, X. Zhu, T. Xiang, and S. Gong, "Towards unsupervised open-set person re-identification," *in Proc. of the IEEE Int. Conf. Image Processing. IEEE*, 2016, pp. 769–773.

[32] Y. Ge, D. Chen, and H. Li, "Mutual mean-teaching: Pseudo label refinery for unsupervised domain adaptation on person re-identification," *arXiv preprint arXiv:2001.01526*, 2020.

[33] M. Jia, Y. Zhai, S. Lu, S. Ma, and J. Zhang, "A similarity inference metric for RGB-infrared cross-modality person re-identification," *arXiv preprint arXiv:2007.01504.*, 2020.

[34] P. Dai, R. Ji, H. Wang, Q. Wu, and Y. Huang, "Cross-modality person re-identification with generative adversarial training." *in Proc. of the Int. Joint Conf. Artificial Intelligence*, vol. 1, 2018, p. 2.

[35] J. Mairal, J. Ponce, G. Sapiro, A. Zisserman, and F. Bach, "Supervised dictionary learning," *Advances in neural information processing systems*, vol. 21, pp. 1033–1040, 2008.

[36] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," *in Proc. of the Int. Conf. Machine Learning*, 2007, pp. 209–216.

[37] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai et al., "Recent advances in convolutional neural networks," *Pattern Recognition*, vol. 77, pp. 354–377, 2018.

[38] Z. Wang, M. Ye, F. Yang, X. Bai, and S. Satoh, "Cascaded sr-gan for scale-adaptive low resolution person re-identification." *in Proc. of the Int. Joint Conf. Artificial Intelligence*, vol. 1, no. 2, 2018, p. 4.

[39] Y.-J. Li, Y.-C. Chen, Y.-Y. Lin, X. Du, and Y.-C. F. Wang, "Recover and identify: A generative dual model for cross-resolution person re-identification," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, 2019, pp. 8090–8099.

[40] Z. Cheng, Q. Dong, S. Gong, and X. Zhu, "Inter-task association critic for cross-resolution person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2020, pp. 2605–2615.

[41] A. Toshev and C. Szegedy, "Deep pose: Human pose estimation via deep neural networks," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2014, pp. 1653–1660.

[42] Y. Sun, X. Wang, and X. Tang, "Deep convolutional network cascade for facial point detection," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2013, pp. 3476–3483.

[43] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," *in Proc. of the European Conf. Comput. Vis.*, 2018, pp. 286–301.

[44] D. J. Finney, *Probit analysis: a statistical treatment of the sigmoid response curve.* Cambridge university press, 1952.

[45] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," *in Proc. of the Int. Conf. Machine Learning*, 2010.

[46] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deep reid: Deep filter pairing neural network for person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2014, pp. 152–159.

[47] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an ensemble of localized features," *in Proc. of the European Conf. Comput. Vis.* Springer, 2008, pp. 262–275.

[48] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino, "Custom pictorial structures for re-identification." *in Proc. of the British Machine Vis. Conf.*, vol. 1, no. 2, 2011, p. 6.

[49] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, 2015, pp. 1116–1124.

[50] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, 2017, pp. 3754–3762.

[51] Y. Sun, L. Zheng, Y. Yang, Q. Tian, and S. Wang, "Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline)," *in Proc. of the European Conf. Comput. Vis.*, 2018, pp. 480–496.

[52] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2017, pp. 4700–4708.

[53] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2016, pp. 770–778.

[54] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2018, pp. 7132–7141.

[55] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, pp. 2672–2680, 2014.

[56] Z. Zhong, L. Zheng, Z. Zheng, S. Li, and Y. Yang, "Camera style adaptation for person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2018, pp. 5157–5166.

[57] N. Martinel, G. Luca Foresti, and C. Micheloni, "Aggregating deep pyramidal representations for person re-identification," *in Proc. of the IEEE Conf. Comput. Vis. and Pattern Recognition*, 2019, pp. 0–0.

[58] Y. Ge, Z. Li, H. Zhao, G. Yin, S. Yi, X. Wang et al., "Fd-gan: Pose-guided feature distilling gan for robust person re-identification," *in Proc. of the Advances in neural information processing systems*, 2018, pp. 1222–1233.

[59] Z. Wang, R. Hu, Y. Yu, J. Jiang, C. Liang, and J. Wang, "Scale-adaptive low-resolution person re-identification via learning a discriminating surface." *in Proc. of the Int. Joint Conf. Artificial Intelligence*, vol. 2, 2016, p. 6.

[60] S. Mao, S. Zhang, and M. Yang, "Resolution-invariant person re-identification," *arXiv preprint arXiv:1906.09748*, 2019.

[61] Y.-J. Li, Y.-C. Chen, Y.-Y. Lin, and Y.-C. F. Wang, "Cross-resolution adversarial dual network for person re-identification and beyond," *arXiv preprint arXiv:2002.09274*, 2020.

[62] K. Han, Y. Huang, Z. Chen, L. Wang, and T. Tan, "Prediction and recovery for adaptive low-resolution person re-identification," *in Proc. of the European Conf. Comput. Vis.* Springer, 2020, pp. 193–209.

[63] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," *in Proc. of the IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2223–2232.

[64] J. Wang, K. Sun, T. Cheng, B. Jiang, C. Deng, Y. Zhao, D. Liu,Y. Mu, M. Tan, X. Wanget al., "Deep high-resolution representation learning for visual recognition," *IEEE Trans. Pattern Anal. and Machine Intelligence*, 2020.

[65] Aguilera, A. Cristhian, J. Francisco, Sappa, D. Angel, Richard, and Toledo, "Learning cross-spectral similarity measures with deep convolutional neural networks," *in Proc. of IEEE Conf. Comput. Vis. and Pattern Recognition*, 2016, pp. 1-9.

[66] J. Gao, C. Xiao, M. Glass, and J. Sun, "COMPOSE: Cross-Modal Pseudo-Siamese Network for Patient Trial Matching," *in Proc. of the 26th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining*, 2020, pp. 803-812.