

# Style Normalization and Restitution for Generalizable Person Re-identification

Xin Jin<sup>1\*</sup> Cuiling Lan<sup>2†</sup> Wenjun Zeng<sup>2</sup> Zhibo Chen<sup>1†</sup> Li Zhang<sup>3</sup>

<sup>1</sup> University of Science and Technology of China <sup>2</sup> Microsoft Research Asia, Beijing, China <sup>3</sup> University of Oxford

jinxustc@mail.ustc.edu.cn {culan,wezeng}@microsoft.com chenzhibo@ustc.edu.cn lz@robots.ox.ac.uk

## Abstract

Existing fully-supervised person re-identification (ReID) methods usually suffer from poor generalization capability caused by domain gaps. The key to solving this problem lies in filtering out **identity-irrelevant** interference and learning domain-invariant person representations. In this paper, we aim to design a generalizable person ReID framework which trains a model on source domains yet is able to generalize/perform well on target domains. To achieve this goal, we propose a simple yet effective **Style Normalization and Restitution (SNR)** module. Specifically, we filter out style variations (e.g., illumination, color contrast) by **Instance Normalization (IN)**. However, such a process inevitably removes discriminative information. We propose to distill identity-relevant feature from the removed information and reconstitute it to the network to ensure high discrimination. For better disentanglement, we enforce a dual causality loss constraint in SNR to encourage the separation of identity-relevant features and identity-irrelevant features. Extensive experiments demonstrate the strong generalization capability of our framework. Our models empowered by the SNR modules significantly outperform the state-of-the-art domain generalization approaches on multiple widely-used person ReID benchmarks, and also show superiority on unsupervised domain adaptation.

## 1. Introduction

Person re-identification (ReID) aims at matching/identifying a specific person across cameras, times, and locations. It facilitates many applications and has attracted a lot of attention.

Abundant approaches have been proposed for supervised person ReID, where a model is trained and tested on different splits of the same dataset [65, 47, 68, 10, 43, 67, 21, 20]. They typically focus on addressing the challenge of geometric misalignment among images caused by diversity of poses/viewpoints. In general, they perform well on the

\*This work was done when Xin Jin was an intern at Microsoft Research Asia.

†Corresponding Author.

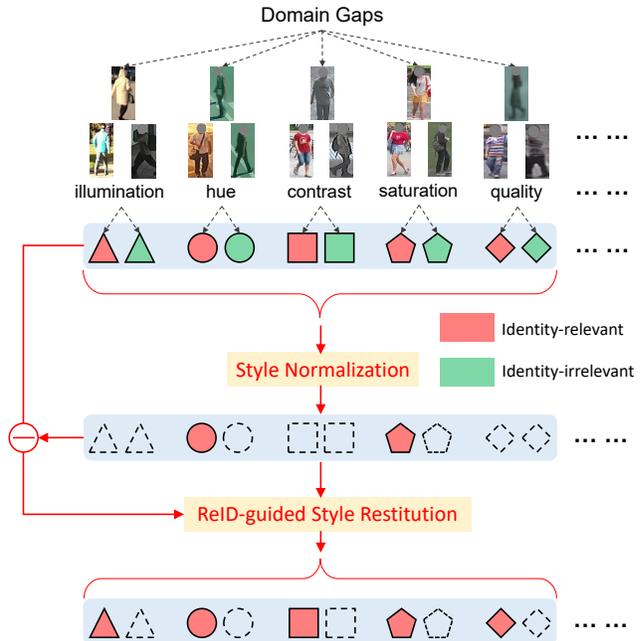


Figure 1: Illustration of motivation and our idea. Person images captured from different cameras and environments present style variations which result in domain gaps. We use style normalization (with Instance Normalization) to alleviate style variations. However, this also results in the loss of some discriminative (identity-relevant) information. We propose to further reconstitute such information from the residual of the original information and the normalized information for generalizable and discriminative person ReID.

trained dataset but suffer from significant performance degradation (poor generalization capability) when testing on a previously unseen dataset. There are usually style discrepancies across domains/datasets which hinder the achievement of high generalization capability. Figure 1 shows some example images<sup>1</sup> from different ReID datasets. The person images are captured by different cameras under different environments (e.g., lighting, seasons). They present a large style discrepancy in terms of illumination, hue, color contrast and saturation, quality/resolution, etc. For a ReID

<sup>1</sup>All faces in the images are masked for anonymization.

system, we expect it to be able to identify the same person even captured in different environments, and distinguish between different people even if their appearance are similar. Both generalization and discrimination capabilities, although seemingly conflicting with each other, are very important for robust ReID.

Considering the existence of domain gaps and poor generalization capability, fully-supervised approaches or settings are not practical for real-world widespread ReID system deployment, where the onsite manual annotation on the target domain data is expensive and hardly feasible. In recent years, some unsupervised domain adaptation (UDA) methods have been studied to adapt a ReID model from source to target domain [53, 50, 35, 42, 7, 64, 60]. UDA models update using *unlabeled* target domain data, emancipating the labelling efforts. However, data collection and model update are still required, adding additional cost.

We mainly focus on the more economical and practical domain generalizable person ReID. Domain generalization (DG) aims to design models that are generalizable to previously unseen domains [40, 19, 45], without having to access the target domain data and labels, and without requiring model updating. Most DG methods assume that the source and target domains have the same label space [22, 26, 40, 44] and they are not applicable to ReID since the target domains for ReID typically have a different label space from the source domains. Generalizable person ReID is challenging which aims to achieve high discrimination capability on *unseen* target domain that may have large domain discrepancy. The study on domain generalizable ReID is rare [45, 19] and remains an open problem. Jia *et al.* [19] and Zhou *et al.* [75] integrate Instance Normalization (IN) in the networks to alleviate the domain discrepancy due to appearance style variations. However, IN inevitably results in the loss of some discriminative features [17, 41], hindering the achievement of high efficiency ReID.

In this paper, we aim to design a generalizable ReID framework which achieves both high generalization capability and discrimination capability. The key is to find a way to disentangle the identity-relevant features and the identity-irrelevant features (*e.g.*, image styles). Figure 1 illustrates our main idea. Considering the domain gaps among image samples, we perform style normalization by means of IN to eliminate style variations. However, the normalization inevitably discards some discriminative information and thus may hamper the ReID performance. From the residual information (which is the difference between the original information and the normalized information), we further distill the identity-relevant information as a compensation to the normalized information. Figure 2 shows our framework with the proposed Style Normalization and Restitution (SNR) modules embedded. To better disentangle the identity-relevant features from the residual, a dual

causality loss constraint is added by ensuring the features after restitution of identity-relevant features to be more discriminative, and the features after compensation of identity-irrelevant features to be less discriminative.

We summarize our main contributions as follows:

- We propose a practical domain generalizable person ReID framework that generalizes well on previously unseen domains/datasets. Particularly, we design a Style Normalization and Restitution (SNR) module. SNR is simple yet effective and can be used as a *plug-and-play* module for existing ReID architectures to enhance their generalization capabilities.
- To facilitate the restitution of identity-relevant features from those discarded in the style normalization phase, we introduce a dual causality loss constraint in SNR for better feature disentanglement.

We validate the effectiveness of the proposed SNR module on multiple widely-used benchmarks and settings. Our models significantly outperform the state-of-the-art domain generalizable person ReID approaches and can also boost the performance of unsupervised domain adaptation for ReID.

## 2. Related Work

**Supervised Person ReID.** In the last decade, fully-supervised person ReID has achieved great progress, especially for deep learning based approaches [47, 25, 68, 10, 43, 67]. These methods usually perform well on the testing set of the source datasets but generalize poorly to previously unseen domains/datasets due to the style discrepancy across domains. This is problematic especially in practical applications, where the target scenes typically have different styles from the source domains and there is no readily available target domain data or annotation for training.

**Unsupervised Domain Adaptation (UDA) for Person ReID.** When the target domain data is accessible, even without annotations, it can be explored for the domain adaptation for enhancing the ReID performance. This requires target domain data collection and model updating. UDA-based ReID methods can be roughly divided into three categories: *style transfer* [5, 56, 35], *attribute recognition* [53, 63, 42], and *target-domain pseudo label estimation* [7, 46, 72, 50, 66, 64]. For pseudo label estimation, recently, Yu *et al.* propose a method called **multilabel reference learning (MAR)** which evaluates the similarity of a pair of images by comparing them to a set of known reference persons to mine hard negative samples [64].

Our proposed domain generalizable SNR module can also be combined with the UDA methods (*e.g.*, by plugging into the UDA backbone) to further enhance the ReID performance. We will demonstrate its effectiveness by combining it with the UDA approach of MAR in Subsection 4.5.

**Domain Generalization (DG).** Domain Generalization is

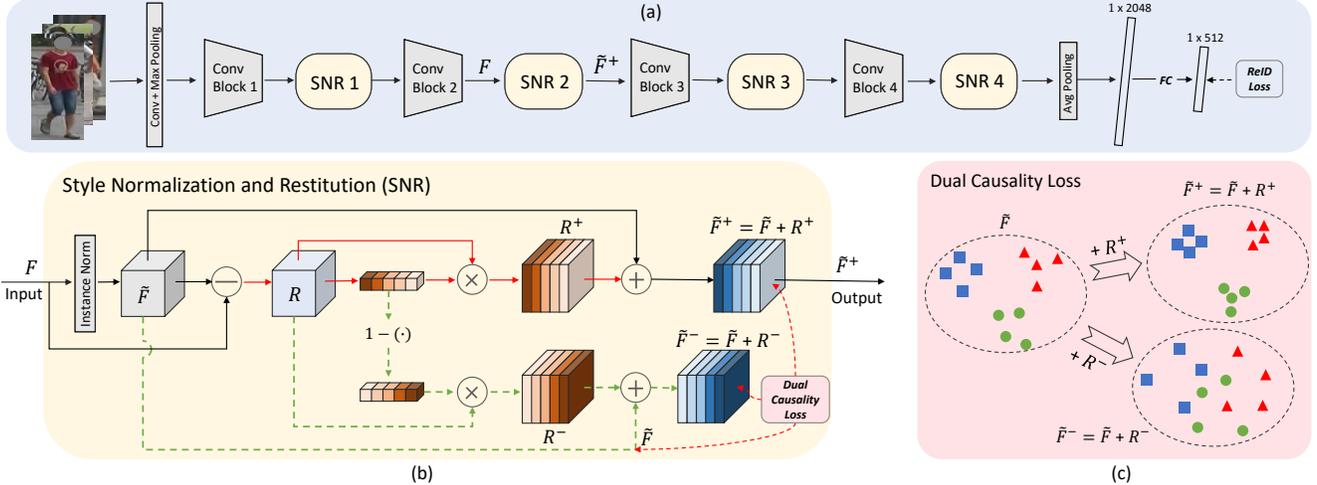


Figure 2: Overall flowchart. (a) Our generalizable person ReID network with the proposed Style Normalization and Restitution (SNR) module being plugged in after some convolutional blocks. Here, we use ResNet-50 as our backbone for illustration. (b) Proposed SNR module. Instance Normalization (IN) is used to eliminate some style discrepancies followed by identity-relevant feature restitution (marked by red solid arrows). Note the branch with dashed green line is only used for enforcing loss constraint and is discarded in inference. (c) Dual causality loss constraint encourages the disentanglement of a residual feature  $R$  to identity-relevant one ( $R^+$ ) and identity-irrelevant one ( $R^-$ ), which enhances and decreases, respectively, the discrimination by adding them to the style normalized feature  $\tilde{F}$ .

a challenging problem of learning models that is generalizable to unseen domains [40, 44]. Muandet *et al.* learn an invariant transformation by minimizing the dissimilarity across source domains [40]. A learning-theoretic analysis shows that reducing dissimilarity improves the generalization ability on new domains. CrossGrad [44] generates pseudo training instances by perturbations in the loss gradients of the domain classifier and category classifier respectively. Most DG methods assume that the source and target domains have the same label space. However, ReID is an open-set problem where the target domains typically have different identities from the source domains, so that the general DG methods could not be directly applied to ReID.

Recently, a strong baseline for domain generalizable person ReID is proposed by simply combing multiple source datasets and training a single CNN [24]. Song *et al.* [45] propose a generalizable person ReID framework by using a meta-learning pipeline to make the model domain invariant. To overcome the inconsistency of label spaces among different datasets, it maintains a training datasets shared memory bank. Instance Normalization (IN) has been widely used in image style transfer [17, 52] and proved that it actually performs a kind of style normalization [41, 17]. Jia *et al.* [19] and Zhou *et al.* [75] apply this idea to ReID to alleviate the domain discrepancy and boost the generalization capability. However, IN inevitably discards some discriminative information. In this paper, we study how to design a generalizable ReID framework that can exploit the merit of IN while avoiding the loss of discriminative information.

### 3. Proposed Generalizable Person ReID

We aim at designing a generalizable and robust person ReID framework. During the training, we have access to one or several annotated source datasets. The trained model will be deployed directly to unseen domains/datasets and is expected to work well with high generalization capability.

Figure 2 shows the overall flowchart of our framework. Particularly, we propose a Style Normalization and Restitution (SNR) module to boost the generalization and discrimination capability of ReID models especially on unseen domains. SNR can be used as a plug-and-play module for existing ReID networks. Taking the widely used ReID network of ResNet-50 [13, 1, 37] as an example (see Figure 2(a)), SNR module is added after each convolutional block. In the SNR module, we first eliminate style discrepancy among samples by Instance Normalization (IN). Then, a dedicated restitution step is proposed to distill identity-relevant (discriminative) features from those previously discarded by IN, and add them to the normalized features. Moreover, for the SNR module, we design a dual causality loss constraint to facilitate the distillation of identity-relevant features from the information discarded by IN.

#### 3.1. Style Normalization and Restitution (SNR)

Person images for ReID could be captured by different cameras under different scenes and environments (*e.g.*, indoor/outdoors, shopping malls, street, sunny/cloudy). As shown in Figure 1, they present style discrepancies (*e.g.*, in illumination, hue, contrast, saturation, quality), especially

for samples from two different datasets/domains. Domain discrepancy between the source and target domain generally hinders the generalization capability of ReID models.

A learning-theoretic analysis shows that reducing dissimilarity improves the generalization ability on new domains [40]. Instance Normalization (IN) performs some kinds of style normalization which reduces the discrepancy/dissimilarity among instances/samples [17, 41], so it can enhance the generalization ability of networks [41, 19, 75]. However, IN inevitably removes some discriminative information and results in weaker discrimination capability [41]. To address this problem, we propose to reconstitute the task-specific discriminative features from the IN removed information, by disentangling it into identity-relevant features and identity-irrelevant features with a dual causality loss constraint (see Figure 2(b)). We elaborate on the designed SNR module hereafter.

For an SNR module, we denote the input (which is a feature map) by  $F \in \mathbb{R}^{h \times w \times c}$  and the output by  $\tilde{F}^+ \in \mathbb{R}^{h \times w \times c}$ , where  $h, w, c$  denote the height, width, and number of channels, respectively.

**Style Normalization Phase.** In SNR, we first try to reduce the domain discrepancy on the input features by performing Instance Normalization [51, 6, 52, 17] as

$$\tilde{F} = \text{IN}(F) = \gamma \left( \frac{F - \mu(F)}{\sigma(F)} \right) + \beta, \quad (1)$$

where  $\mu(\cdot)$  and  $\sigma(\cdot)$  denote the mean and standard deviation computed across spatial dimensions independently for each channel and each *sample/instance*,  $\gamma, \beta \in \mathbb{R}^c$  are parameters learned from data. IN could filter out some instance-specific style information from the content. With IN taking place in the feature space, Huang *et al.* [17] have argued and experimentally shown that IN has more profound impacts than a simple contrast normalization and it performs a form of *style normalization* by normalizing feature statistics.

**Style Restitution Phase.** IN reduces style discrepancy and boosts the generalization capability. However, with the mathematical operations being deterministic and task-irrelevant, it inevitably discards some discriminative (task-relevant) information for ReID. We propose to reconstitute the identity-relevant feature to the network by distilling it from the residual feature  $R$ .  $R$  is defined as

$$R = F - \tilde{F}, \quad (2)$$

which denotes the difference between the original input feature  $F$  and the style normalized feature  $\tilde{F}$ .

Given  $R$ , we further disentangle it into two parts: identity-relevant feature  $R^+ \in \mathbb{R}^{h \times w \times c}$  and identity-irrelevant feature  $R^- \in \mathbb{R}^{h \times w \times c}$  through masking  $R$  by a learned channel attention vector  $\mathbf{a} = [a_1, a_2, \dots, a_c] \in \mathbb{R}^c$ :

$$\begin{aligned} R^+(\cdot, \cdot, k) &= a_k R(\cdot, \cdot, k), \\ R^-(\cdot, \cdot, k) &= (1 - a_k) R(\cdot, \cdot, k), \end{aligned} \quad (3)$$

where  $R(\cdot, \cdot, k) \in \mathbb{R}^{h \times w}$  denotes the  $k^{\text{th}}$  channel of feature map  $R$ ,  $k = 1, 2, \dots, c$ . We expect the channel attention vector  $\mathbf{a}$  to enable the adaptive distillation of the identity-relevant features for restitution, and derive it by SE-like [16] channel attention as

$$\mathbf{a} = g(R) = \sigma(W_2 \delta(W_1 \text{pool}(R))), \quad (4)$$

which consists of a global average pooling layer followed by two FC layers that are parameterized by  $W_2 \in \mathbb{R}^{(c/r) \times c}$  and  $W_1 \in \mathbb{R}^{c \times (c/r)}$  which are followed by ReLU activation function  $\delta(\cdot)$  and sigmoid activation function  $\sigma(\cdot)$ , respectively. To reduce the number of parameters, a dimension reduction ratio  $r$  is used and is set to 16.

By adding the distilled identity-relevant feature  $R^+$  to the style normalized feature  $\tilde{F}$ , we obtain the output feature  $\tilde{F}^+$  of the SNR module as

$$\tilde{F}^+ = \tilde{F} + R^+. \quad (5)$$

**Dual Causality Loss Constraint.** In order to facilitate the disentanglement of identity-relevant feature and identity-irrelevant feature, we design a dual causality loss constraint by comparing the discrimination capability of features *before* and *after* the restitution. As illustrated in Figure 2(c), the main idea is that: after restituting the identity-relevant feature  $R^+$  to the normalized feature  $\tilde{F}$ , the feature becomes more discriminative; On the other hand, after restituting the identity-irrelevant feature  $R^-$  to the normalized feature  $\tilde{F}$ , the feature should become less discriminative. We achieve this by defining a dual causality loss  $\mathcal{L}_{SNR}$  which consists of *clarification loss*  $\mathcal{L}_{SNR}^+$  and *destruction loss*  $\mathcal{L}_{SNR}^-$ , *i.e.*,  $\mathcal{L}_{SNR} = \mathcal{L}_{SNR}^+ + \mathcal{L}_{SNR}^-$ .

Within a mini-batch, we sample three images, *i.e.*, an anchor sample  $a$ , a positive sample  $p$  that has the same identity as the anchor sample, and a negative sample  $n$  that has a different identity from the anchor sample. For simplicity, we differentiate the three samples by subscript. For example, the style normalized feature of sample  $a$  is denoted by  $\tilde{F}_a$ .

Intuitively, adding the identity-relevant feature  $R^+$  to the normalized feature  $\tilde{F}$ , which we refer to as *enhanced feature*  $\tilde{F}^+ = \tilde{F} + R^+$ , results in better discrimination capability — the sample features with same identities are closer and those with different identities are farther apart. We calculate the distances between samples on a spatially average pooled feature to avoid the distraction caused by spatial misalignment among samples (*e.g.*, due to different poses/viewpoints). We denote the spatially average pooled feature of  $\tilde{F}$  and  $\tilde{F}^+$  as  $\tilde{\mathbf{f}} = \text{pool}(\tilde{F})$ ,  $\tilde{\mathbf{f}}^+ = \text{pool}(\tilde{F}^+)$ , respectively. The *clarification loss* is thus defined as

$$\begin{aligned} \mathcal{L}_{SNR}^+ &= \text{Softplus}(d(\tilde{\mathbf{f}}_a^+, \tilde{\mathbf{f}}_p^+) - d(\tilde{\mathbf{f}}_a, \tilde{\mathbf{f}}_p)) \\ &+ \text{Softplus}(d(\tilde{\mathbf{f}}_a, \tilde{\mathbf{f}}_n) - d(\tilde{\mathbf{f}}_a^+, \tilde{\mathbf{f}}_n^+)), \end{aligned} \quad (6)$$

where  $d(\mathbf{x}, \mathbf{y})$  denotes the distance between  $\mathbf{x}$  and  $\mathbf{y}$  which is defined as  $d(\mathbf{x}, \mathbf{y}) = 0.5 - \mathbf{x}^T \mathbf{y} / (2\|\mathbf{x}\| \|\mathbf{y}\|)$ .  $Softplus(\cdot) = \ln(1 + \exp(\cdot))$  is a monotonically increasing function that aims to reduce the optimization difficulty by avoiding negative loss values.

On the other hand, we expect that the adding of the identity-irrelevant feature  $R^-$  to the normalized feature  $\tilde{F}$ , which we refer to as *contaminated feature*  $\tilde{F}^- = \tilde{F} + R^-$ , could decrease the discrimination capability. In comparison with the normalized feature  $\tilde{F}$  before the compensation, we expect that adding  $R^-$  would push the sample features with same identities farther apart and pull those with different identities closer. We denote the spatially average pooled feature of  $\tilde{F}^-$  as  $\tilde{\mathbf{f}}^- = pool(\tilde{F}^-)$ . The *destruction loss* is:

$$\begin{aligned} \mathcal{L}_{SNR}^- = & Softplus(d(\tilde{\mathbf{f}}_a, \tilde{\mathbf{f}}_p) - d(\tilde{\mathbf{f}}_a^-, \tilde{\mathbf{f}}_p^-)) \\ & + Softplus(d(\tilde{\mathbf{f}}_a^-, \tilde{\mathbf{f}}_n^-) - d(\tilde{\mathbf{f}}_a, \tilde{\mathbf{f}}_n)). \end{aligned} \quad (7)$$

### 3.2. Joint Training

We use the commonly used ResNet-50 as a base ReID network and insert the proposed SNR module after each convolution block (in total four convolution blocks/stages)(see Figure 2(a)). We train the entire network in an end-to-end manner. The overall loss is

$$\mathcal{L} = \mathcal{L}_{ReID} + \sum_{b=1}^4 \lambda_b \mathcal{L}_{SNR}^b, \quad (8)$$

where  $\mathcal{L}_{SNR}^b$  denotes the dual causality loss for the  $b^{th}$  SNR module.  $\mathcal{L}_{ReID}$  denotes the widely-used ReID Loss (classification loss [48, 9], and triplet loss with batch hard mining [14]) on the ReID feature vectors.  $\lambda_b$  is a weight which controls the relative importance of the regularization at stage  $b$ . In considering that the features of stage 3 and 4 are more relevant to the task (high-level semantics), we experimentally set  $\lambda_3, \lambda_4$  to 0.5, and  $\lambda_1, \lambda_2$  to 0.1.

## 4. Experiments

In this section, we first describe the datasets and evaluation metrics in Subsection 4.1. Then, for generalizable ReID, we validate the effectiveness of SNR in Subsection 4.2 and study its design choices in Subsection 4.3. We conduct visualization analysis in Subsection 4.4. Subsection 4.5 shows the comparisons of our schemes with the state-of-the-art approaches for both generalizable person ReID and unsupervised domain adaption ReID, respectively. In Subsection 4.6, we further validate the effectiveness of applying the SNR modules to another backbone network and to cross modality (Infrared-RGB) person ReID.

We use ResNet-50 [13, 1, 67, 37] as our base network for both baselines and our schemes. We build a strong baseline *Baseline* with some commonly used tricks integrated.

### 4.1. Datasets and Evaluation Metrics

To evaluate the generalization ability of our approach and to be consistent with what were done in prior works for performance comparisons, we conduct extensive experiments on commonly used public ReID datasets, including Market1501 [69], DukeMTMC-reID [71], CUHK03 [28], the large-scale MSMT17 [56], and four small-scale ReID datasets of PRID [15], GRID [36], VIPeR [11], and i-LIDS [57]. We denote Market1501 by M, DukeMTMC-reID by Duke or D, and CUHK03 by C for simplicity.

We follow common practices and use the cumulative matching characteristics (CMC) at Rank-1, and mean average precision (mAP) to evaluate the performance.

### 4.2. Ablation Study

We perform comprehensive ablation studies to demonstrate the effectiveness of the SNR module and its dual causality loss constraint. We mimic the real-world scenario for generalizable person ReID, where a model is trained on some source dataset(s) A while tested on previously unseen dataset B. We denote this as A→B. We have several experimental settings to evaluate the generalization capability, e.g., Market1501→Duke and others, Duke→Market1501 and others, M+D+C+MSMT17→others. Our settings cover both single source dataset for training and multiple source datasets for training.

**Effectiveness of Our SNR.** Here we compare several schemes. **Baseline:** a strong baseline based on ResNet-50. **Baseline-A-IN:** a naive model where we replace all the Batch Normalization(BN) [18] layers in *Baseline* by Instance Normalization(IN). **Baseline-IBN:** Similar to IBN-Net (IBN-b) [41] and OSNet [75], we add IN only to the last layers of Conv1 and Conv2 blocks of *Baseline* respectively. **Baseline-A-SN:** a model where we replace all the BN layers in *Baseline* by Switchable Normalization (SN). SN [38] can be regarded as an adaptive ensemble version of normalization techniques of IN, BN, and LN (Layer Normalization) [2]. **Baseline-IN:** four IN layers are added after the first four convolutional blocks/stages of *Baseline* respectively. **Baseline-SNR:** our final scheme where four SNR modules are added after the first four convolutional blocks/stages of *Baseline* respectively (see Figure 2(a)). We also refer to it as **SNR** for simplicity. Table 5 shows the results. We have the following observations/conclusions:

1) *Baseline-A-IN* improves *Baseline* by **4.3%** in mAP for Market1501→Duke, and **4.7%** in mAP for Duke→Market1501. Other IN-related baselines also bring gains, which demonstrates the effectiveness of IN for improving the generalization capability for ReID. But, IN also inevitably discards some discriminative (identity-relevant) information and we can see it clearly decreases the performance of *Baseline-A-IN*, *Baseline-IBN* and *Baseline-IN* for the same-domain ReID (e.g., Market1501→Market1501).

Table 1: Performance (%) comparisons of our scheme and others to demonstrate the effectiveness of our SNR module for generalizable person ReID. The rows denote source dataset(s) for training and the columns correspond to different target datasets for testing. We mask the results of supervised ReID by gray where the testing domain has been seen in training. Due to space limitation, we only show a portion of the results here and more comparisons can be found in **Supplementary**.

Source	Method	Target: Market1501		Target: Duke		Target: PRID		Target: GRID		Target: VIPeR		Target: iLIDs	
		mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1
Market1501 (M)	Baseline	82.8	93.2	19.8	35.3	13.7	6.0	25.8	16.0	37.6	28.5	61.5	53.3
	Baseline-A-IN	75.3	89.8	24.1	42.7	33.9	21.0	35.6	27.2	38.1	29.1	64.2	55.0
	Baseline-IBN	81.1	92.2	21.5	39.2	19.1	12.0	27.5	19.2	32.1	23.4	58.3	48.3
	Baseline-A-SN	83.2	93.9	20.1	38.0	35.4	25.0	29.0	22.0	32.2	23.4	53.4	43.3
	Baseline-IN	79.5	90.9	25.1	44.9	35.0	25.0	35.7	27.8	35.1	27.5	64.0	54.2
	<b>Baseline-SNR (Ours)</b>	<b>84.7</b>	<b>94.4</b>	<b>33.6</b>	<b>55.1</b>	<b>42.2</b>	<b>30.0</b>	<b>36.7</b>	<b>29.0</b>	<b>42.3</b>	<b>32.3</b>	<b>65.6</b>	<b>56.7</b>
Duke (D)	Baseline	21.8	48.3	71.2	83.4	15.7	11.0	14.5	8.8	37.0	26.9	68.3	58.3
	Baseline-A-IN	26.5	56.0	64.5	78.9	38.6	29.0	19.6	13.6	35.1	27.2	67.4	56.7
	Baseline-IBN	24.6	52.5	69.5	81.4	27.4	19.0	19.9	12.0	32.8	23.4	63.5	61.7
	Baseline-A-SN	25.3	55.0	73.0	85.9	41.4	32.0	18.8	12.8	31.3	24.1	64.8	63.3
	Baseline-IN	27.2	58.5	68.9	80.4	40.5	27.0	20.3	13.2	34.6	26.3	70.6	65.0
	<b>Baseline-SNR (Ours)</b>	<b>33.9</b>	<b>66.7</b>	<b>72.9</b>	<b>84.4</b>	<b>45.4</b>	<b>35.0</b>	<b>35.3</b>	<b>26.0</b>	<b>41.2</b>	<b>32.6</b>	<b>79.3</b>	<b>68.7</b>
M + D + CUHK03 + MSMT17	Baseline	72.4	88.7	70.1	83.8	39.0	28.0	29.6	20.8	52.1	41.5	89.0	85.0
	<b>Baseline-SNR (Ours)</b>	<b>82.3</b>	<b>93.4</b>	<b>73.2</b>	<b>85.5</b>	<b>60.0</b>	<b>49.0</b>	<b>41.3</b>	<b>30.4</b>	<b>65.0</b>	<b>55.1</b>	<b>91.9</b>	<b>87.0</b>

*Baseline-A-SN* learns the combination weights of IN, BN, and LN in the training dataset and thus has superior performance in the same domain, but it does not have dedicated design for boosting the generalization capability.

2) Thanks to the compensation of the identity-relevant information through the proposed *restitution step*, our final scheme *Baseline-SNR* achieves superior generalization capability, which significantly outperforms all the baseline schemes. In particular, *Baseline-SNR* outperforms *Baseline-IN* by **8.5%**, **6.7%**, **15.0%** in mAP for M→D, D→M, and D→GRID, respectively.

3) The generalization performance on previously unseen target domain increases consistently as the number of source datasets increases. When all the four source datasets are used (the large-scale MSMT17 [56] also included), we have a very strong baseline (*i.e.*, 52.1% in mAP on VIPeR dataset vs. 37.6% when Market1501 alone is used as source). Interestingly, our method still significantly outperforms the strong baseline *Baseline*, even by **21.0%** in mAP on PRID dataset, demonstrating SNR’s effectiveness.

4) The performance of different schemes with respects to PRID/GRID varies greatly and the mAPs are all relatively low, which is caused by the large style discrepancy between PRID/GRID and other datasets. For such challenging cases, our scheme still outperforms *Baseline-IN* significantly by **7.2%** and **4.9%** in mAP for M→PRID and D→PRID, respectively.

5) For supervised ReID (masked by gray), our scheme also clearly outperforms *Baseline* by **1.9%** and **1.7%** in mAP for M→M and D→D, respectively. That is because there is also style discrepancy within the source domain.

**Influence of Dual Causality Loss Constraint.** We study the effectiveness of the proposed dual causality loss  $\mathcal{L}_{SNR}$  which consists of *clarification loss*  $\mathcal{L}_{SNR}^+$  and *destruction loss*  $\mathcal{L}_{SNR}^-$ . Table 2a shows the results. Our final scheme *SNR* with the dual causality loss  $\mathcal{L}_{SNR}$  outperforms that without such constraints (*i.e.*, scheme *SNR w/o*  $\mathcal{L}_{SNR}$ ) by

**7.5%** and **4.7%** in mAP for M→D and D→M, respectively. Such constraints facilitate the disentanglement of identity-relevant/identity-irrelevant features. In addition, both the clarification loss  $\mathcal{L}_{SNR}^+$  and the destruction loss  $\mathcal{L}_{SNR}^-$  are vital to SNR and they are complementary and jointly contribute to a superior performance.

**Complexity.** The model size of our final scheme *SNR* is very similar to that of *Baseline* (24.74 M vs. 24.56 M).

### 4.3. Design Choices of SNR

**Which Stage to Add SNR?** We compare the cases of adding a single SNR module to a different convolutional block/stage, and to all the four stages (*i.e.*, stage-1 ~ 4) of the ResNet-50 (see Figure 2(a)). The module is added after the last layer of a convolutional block/stage. As Table 2b shows, in comparison with *Baseline*, the improvement from adding SNR is significant on stage-3 and stage-4 and is a little smaller on stage-1 and stage-2. When SNR is added to all the four stages, we achieve the best performance.

**Influence of Disentanglement Design.** In our SNR module, as described in (3)(4) of Subsection 3.1, we use  $g(\cdot)$ , and its complementary one  $1 - g(\cdot)$  as masks to extract identity-relevant feature  $R^+$  and identity-irrelevant feature  $R^-$  from the residual feature  $R$ . Here, we study the influence of different disentanglement designs within SNR. ***SNR<sub>conv</sub>***: we disentangle the residual feature  $R$  through  $1 \times 1$  convolutional layer followed by non-linear ReLU activation, *i.e.*,  $R^+ = ReLU(W^+R)$ ,  $R^- = ReLU(W^-R)$ . ***SNR<sub>g(\cdot)^2</sub>***: we use two unshared gates  $g(\cdot)^+$ ,  $g(\cdot)^-$  to obtain  $R^+$  and  $R^-$  respectively. Table 2c shows the results. We observe that (1) ours outperforms *SNR<sub>conv</sub>* by **3.9%** and **4.5%** in mAP for M→D and D→M, respectively, demonstrating the benefit of content-adaptive design; (2) ours outperforms *SNR<sub>g(\cdot)^2</sub>* by **2.4%**/**2.9%** in mAP on the unseen target Duke/Market1501, demonstrating the benefit of the design which encourages interaction between  $R^+$  and  $R^-$ .

Table 2: Effectiveness of dual causality loss constraint (a), and study on design choices of SNR (b) and (c).

(a) Study on the dual causality loss constraint.					(b) Study on which stage to add SNR.					(c) Disentanglement designs in SNR.				
Method	M→D		D→M		Method	M→D		D→M		Method	M→D		D→M	
	mAP	Rank-1	mAP	Rank-1		mAP	Rank-1	mAP	Rank-1		mAP	Rank-1	mAP	Rank-1
Baseline	19.8	35.3	21.8	48.3	Baseline	19.8	35.3	21.8	48.3	Baseline	19.8	35.3	21.8	48.3
SNR w/o $\mathcal{L}_{SNR}$	26.1	45.0	29.2	57.4	stage-1	23.7	42.8	27.6	57.7	SNR <sub>conv</sub>	29.7	51.1	29.4	61.7
SNR w/o $\mathcal{L}_{SNR}^+$	28.8	48.9	30.2	59.8	stage-2	24.0	44.4	28.6	58.8	SNR <sub>g(\cdot)^2</sub>	31.2	52.9	31.0	63.8
SNR w/o $\mathcal{L}_{SNR}^-$	28.0	48.1	30.3	59.1	stage-3	26.4	46.3	29.5	60.7	<b>SNR</b>	<b>33.6</b>	<b>55.1</b>	<b>33.9</b>	<b>66.7</b>
<b>SNR</b>	<b>33.6</b>	<b>55.1</b>	<b>33.9</b>	<b>66.7</b>	stage-4	26.2	45.8	29.4	59.7					
					<b>stages-all</b>	<b>33.6</b>	<b>55.1</b>	<b>33.9</b>	<b>66.7</b>					

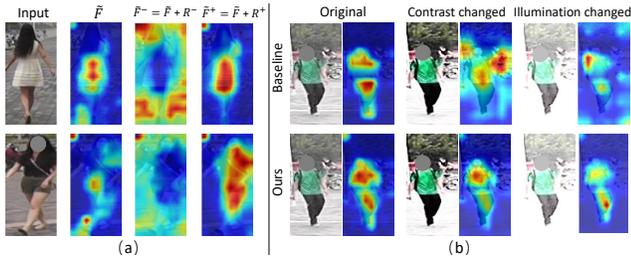


Figure 3: (a) Activation maps of different features within an SNR module (SNR 3). They show SNR can disentangle the identity-relevant/irrelevant features well. (b) Activation maps of our scheme (bottom) and the strong baseline *Baseline* (top) corresponding to images of varied styles. Our maps are more consistent/invariant to style variants.

#### 4.4. Visualization

**Feature Map Visualization.** To better understand how an SNR module works, we visualize the intermediate feature maps of the third SNR module (SNR 3). Following [75, 70], we get each activation map by summarizing the feature maps along channels followed by a spatial  $\ell_2$  normalization.

Figure 6(a) shows the activation maps of normalized feature  $\tilde{F}$ , enhanced feature  $\tilde{F}^+ = \tilde{F} + R^+$ , and contaminated feature  $\tilde{F}^- = \tilde{F} + R^-$ , respectively. We see that after adding the identity-irrelevant feature  $R^-$ , the contaminated feature  $\tilde{F}^-$  has high response mainly on background. In contrast, the enhanced feature  $\tilde{F}^+$  with the restitution of identity-relevant feature  $R^+$  has high responses on regions of the human body, better capturing discriminative regions.

Moreover, in Figure 6(b), we further compare the activation maps  $\tilde{F}^+$  of our scheme and those of the strong baseline *Baseline* by varying the styles of input images (e.g., contrast, illumination, saturation). We can see that, for the images with different styles, the activation maps of our scheme are more consistent/invariant than those of *Baseline*. In contrast, the activation maps of *Baseline* are more disorganized and are easily affected by style variants. These indicate our scheme is more robust to style variations.

**Visualization of Feature Distributions.** In Figure 4, we visualize the distribution of the features from the 3<sup>rd</sup> SNR module of our network using t-SNE [39]. They denote the distributions of features for (a) input  $F$ , (b) style normalized feature  $\tilde{F}$ , and (c) output  $\tilde{F}^+$  of the SNR module. We observe that, (a) before SNR, the extracted features from

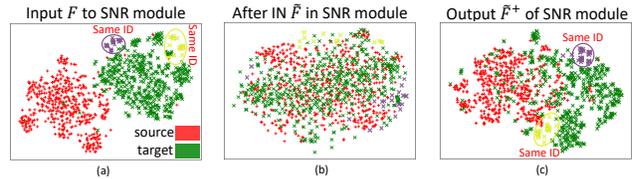


Figure 4: Visualization of distributions of intermediate features before/within/after the SNR module using the tool of t-SNE [39]. ‘Red’/‘green’ nodes: samples from source dataset Market1501/unseen target dataset Duke.

two datasets (‘red’: source training dataset Market1501; ‘green’: unseen target dataset Duke) are largely separately distributed and have an obvious *domain gap*. (b) Within the SNR module, after IN, this *domain gap* has been eliminated. But the samples of the same identity (‘yellow’ and ‘purple’ colored nodes denote two identities respectively) become dispersive. (c) After the restitution of identity-relevant features, not only has the domain gap of feature distributions been shrunk, but also the feature distribution of samples with same identity become more compact than that in (b).

#### 4.5. Comparison with State-of-the-Arts

Thanks to the capability of reducing style discrepancy and restitution of identity-relevant features, our proposed SNR module can enhance the generalization ability and maintain the discriminative ability of ReID networks. It can be used for generalizable person ReID, *i.e.*, domain generalization (DG), and can also be used to build the backbone networks for unsupervised domain adaptation (UDA) for person ReID. We evaluate the effectiveness of SNR on both DG-ReID and UDA-ReID by comparing with the state-of-the-art approaches in Table 6.

**Domain generalizable person ReID** is very attractive in practical applications, which supports “train once and run everywhere”. However, there are very few works in this field [45, 19, 75, 24]. Thanks to the exploration of the style normalization and restitution, our scheme *SNR(Ours)* significantly outperforms the second best method *OSNet-IBN* [75] by **6.9%** and **7.8%** for Market1501→Duke and Duke→Market1501 in mAP, respectively. *OSNet-IBN* adds Instance Normalization (IN) to the lower layers of their proposed OSNet following [41]. However, this does not overcome the intrinsic shortcoming of IN and is not optimal.

Song *et al.* [45] also explore domain generalizable

Table 3: Performance (%) comparisons with the state-of-the-art approaches for the Domain Generalizable Person ReID (top rows) and the Unsupervised Domain Adaptation Person ReID (bottom rows), respectively. “(U)” denotes “unlabeled”. We mask the schemes that use our *Baseline* and those that use our SNR modules by gray, which provides fair comparison.

	Method	Venue	Source	Target: Duke		Source	Target: Market1501		
				mAP	Rank-1		mAP	Rank-1	
Domain Generalization (w/o using target data)	OSNet-IBN [75]	ICCV'19	Market1501	26.7	48.5	Duke	26.1	57.7	
	Baseline	This work	Market1501	19.8	35.3	Duke	21.8	48.3	
	Baseline-IBN [19]	BMVC'19	Market1501	21.5	39.2	Duke	24.6	52.5	
	<b>SNR(Ours)</b>	This work	Market1501	<b>33.6</b>	<b>55.1</b>	Duke	<b>33.9</b>	<b>66.7</b>	
	StrongBaseline [24]	ArXiv'19	MSMT17	43.3	64.5	MSMT17	36.6	64.8	
	OSNet-IBN [75]	ICCV'19	MSMT17	45.6	67.4	MSMT17	37.2	66.5	
	Baseline	This work	MSMT17	39.1	60.4	MSMT17	33.8	59.9	
	<b>SNR(Ours)</b>	This work	MSMT17	<b>50.0</b>	<b>69.2</b>	MSMT17	<b>41.4</b>	<b>70.1</b>	
	Unsupervised Domain Adaptation (using unlabeled target data)	ATNet [35]	CVPR'19	Market1501 + Duke (U)	24.9	45.1	Duke + Market1501 (U)	25.6	55.7
		CamStyle [74]	TIP'19	Market1501 + Duke (U)	25.1	48.4	Duke + Market1501 (U)	27.4	58.8
ARN [30]		CVPRW'19	Market1501 + Duke (U)	33.4	60.2	Duke + Market1501 (U)	39.4	70.3	
ECN [73]		CVPR'19	Market1501 + Duke (U)	40.4	63.3	Duke + Market1501 (U)	43.0	75.1	
PAST [66]		ICCV'19	Market1501 + Duke (U)	54.3	72.4	Duke + Market1501 (U)	54.6	78.4	
SSG [8]		ICCV'19	Market1501 + Duke (U)	53.4	73.0	Duke + Market1501 (U)	58.3	80.0	
Baseline+MAR [64]		This work	Market1501 + Duke (U)	35.2	56.5	Duke + Market1501 (U)	37.2	62.4	
<b>SNR(Ours)+MAR [64]</b>		This work	Market1501 + Duke (U)	<b>58.1</b>	<b>76.3</b>	Duke + Market1501 (U)	<b>61.7</b>	<b>82.8</b>	
MAR [64]		CVPR'19	MSMT17 + Duke (U)	48.0	67.1	MSMT17 + Market1501 (U)	40.0	67.7	
PAUL [60]		CVPR'19	MSMT17 + Duke (U)	53.2	72.0	MSMT17 + Market1501 (U)	40.1	68.5	
Baseline+MAR [64]		This work	MSMT17 + Duke (U)	46.2	66.3	MSMT17 + Market1501 (U)	39.4	66.9	
<b>SNR(Ours) + MAR [64]</b>		This work	MSMT17 + Duke (U)	<b>61.6</b>	<b>78.2</b>	MSMT17 + Market1501 (U)	<b>65.9</b>	<b>85.5</b>	

person ReID and propose a Domain-Invariant Mapping Network (*DIMN*) to learn the mapping between a person image and its identity classifier with a meta-learning pipeline. We follow [45] and train *SNR* on the same five datasets (M+D+C+CUHK02[27]+CUHK-SYSU[59]). *SNR* outperforms *DIMN* by **14.6%/6.6%/1.2%/11.5%** in mAP and **12.9%/10.9%/1.7%/13.9%** in Rank-1 on the PRID/GRID/VIPeR/i-LIDS.

**Unsupervised domain adaptation for ReID** has been extensively studied where the unlabeled target data is also used for training. We follow the most commonly-used source→target setting [73, 35, 75, 64, 60] for comparison. We take *SNR* (see Figure 2(a)) as the backbone followed by a domain adaptation strategy MAR [64] for domain adaptation, which we denote as *SNR(Ours)+MAR* [64]. For comparison, we take our strong *Baseline* as the backbone followed by MAR, which we denote as *Baseline+MAR*, to evaluate the effectiveness of the proposed SNR modules. We can see that *SNR(Ours)+MAR* [64] significantly outperforms the second-best UDA ReID method by **3.8%, 3.4%** in mAP for Market1501+Duke(U)→Duke and Duke+Market1501(U)→Market1501, respectively. In addition, *SNR(Ours)+MAR* outperforms *Baseline+MAR* by **22.9%, 24.5%** in mAP. Similar trends can be found for MSMT17+Duke(U)→Duke and MSMT17+Market1501(U)→Market1501.

In general, as a plug-and-play module, SNR clearly enhances the generalization capability of ReID networks.

#### 4.6. Extension

**Performance on Other Backbone.** We add SNR into the recently proposed lightweight ReID network OSNet [75] and observe that by simply inserting SNR modules between

the OS-Blocks, the new scheme *OSNet-SNR* outperforms their model *OSNet-IBN* by **5.0%** and **5.5%** in mAP for M→D and D→M, respectively (see **Supplementary**).

**RGB-Infrared Cross-Modality Person ReID.** To further demonstrate the capability of SNR in handling images with large style variations, we conduct experiment on a more challenging RGB-Infrared cross-modality person ReID task on benchmark dataset SYSU-MM01 [58]. Our scheme which integrates SNR to *Baseline* outperforms *Baseline* significantly by **8.4%, 8.2%, 11.0%**, and **11.5%** in mAP under 4 different settings, and also achieves the state-of-the-art performance (see **Supplementary** for more details).

## 5. Conclusion

In this paper, we propose a generalizable person ReID framework to enable effective ReID. A Style Normalization and Restitution (SNR) module is introduced to exploit the merit of Instance Normalization (IN) that filters out the interference from style variations, and reconstitute the identity-relevant features that are discarded by IN. To efficiently disentangle the identity-relevant and -irrelevant features, we further design a dual causality loss constraint in SNR. Extensive experiments on several benchmarks/settings demonstrate the effectiveness of SNR. Our framework with SNR embedded achieves the best performance on both domain generalization and unsupervised domain adaptation ReID. Moreover, we have also verified SNR’s effectiveness on RGB-Infrared ReID task, and on another backbone.

## 6. Acknowledgments

This work was supported in part by NSFC under Grant U1908209, 61632001 and the National Key Research and Development Program of China 2018AAA0101400.

## References

- [1] Jon Almazan, Bojana Gajic, Naila Murray, and Diane Larlus. Re-id done right: towards good practices for person re-identification. *arXiv preprint arXiv:1801.05339*, 2018.
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- [3] Pingyang Dai, Rongrong Ji, Haibin Wang, Qiong Wu, and Yuyu Huang. Cross-modality person re-identification with generative adversarial training. In *IJCAI*, pages 677–683, 2018.
- [4] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. 2005.
- [5] Weijian Deng, Liang Zheng, Qixiang Ye, Guoliang Kang, Yi Yang, and Jianbin Jiao. Image-image domain adaptation with preserved self-similarity and domain-dissimilarity for person re-identification. In *CVPR*, 2018.
- [6] Vincent Dumoulin, Jonathon Shlens, and Manjunath Kudlur. A learned representation for artistic style. *ICLR*, 2017.
- [7] Hehe Fan, Liang Zheng, Chenggang Yan, and Yi Yang. Unsupervised person re-identification: Clustering and fine-tuning. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 2018.
- [8] Yang Fu, Yunchao Wei, Guanshuo Wang, Xi Zhou, Honghui Shi, and Thomas S. Huang. Self-similarity grouping: A simple unsupervised cross domain adaptation approach for person re-identification. *ICCV*, abs/1811.10144, 2019.
- [9] Yang Fu, Yunchao Wei, Yuqian Zhou, et al. Horizontal pyramid matching for person re-identification. In *AAAI*, 2019.
- [10] Yixiao Ge, Zhuowan Li, Haiyu Zhao, et al. Fd-gan: Pose-guided feature distilling gan for robust person re-identification. In *NeurIPS*, 2018.
- [11] Douglas Gray and Hai Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *ECCV*, pages 262–275. Springer, 2008.
- [12] Yi Hao, Nannan Wang, Jie Li, and Xinbo Gao. Hsme: Hypersphere manifold embedding for visible thermal person re-identification. In *AAAI*, volume 33, pages 8385–8392, 2019.
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, et al. Deep residual learning for image recognition. In *CVPR*, 2016.
- [14] Alexander Hermans, Lucas Beyler, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.
- [15] Martin Hirzer, Csaba Belezna, Peter M Roth, and Horst Bischof. Person re-identification by descriptive and discriminative classification. In *SCIA*, pages 91–102. Springer, 2011.
- [16] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, pages 7132–7141, 2018.
- [17] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *ICCV*, pages 1501–1510, 2017.
- [18] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, 2015.
- [19] Jieru Jia, Qiuqi Ruan, and Timothy M Hospedales. Frustratingly easy person re-identification: Generalizing person re-id in practice. *BMVC*, 2019.
- [20] Xin Jin, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Uncertainty-aware multi-shot knowledge distillation for image-based object re-identification. In *AAAI*, 2020.
- [21] Xin Jin, Cuiling Lan, Wenjun Zeng, Guoqiang Wei, and Zhibo Chen. Semantics-aligned representation learning for person re-identification. In *AAAI*, 2020.
- [22] Aditya Khosla, Tinghui Zhou, Tomasz Malisiewicz, Alexei A Efros, and Antonio Torralba. Undoing the damage of dataset bias. In *ECCV*, pages 158–171. Springer, 2012.
- [23] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.
- [24] Devinder Kumar, Parthipan Siva, Paul Marchwica, and Alexander Wong. Fairest of them all: Establishing a strong baseline for cross-domain person reid. *arXiv preprint arXiv:1907.12016*, 2019.
- [25] Dangwei Li, Xiaotang Chen, Zhang Zhang, et al. Learning deep context-aware features over body and latent parts for person re-identification. In *CVPR*, 2017.
- [26] Da Li, Yongxin Yang, Yi-Zhe Song, and Timothy M Hospedales. Learning to generalize: Meta-learning for domain generalization. In *AAAI*, 2018.
- [27] Wei Li and Xiaogang Wang. Locally aligned feature transforms across views. In *CVPR*, pages 3594–3601, 2013.
- [28] Wei Li, Rui Zhao, Lu Tian, et al. Deepreid: Deep filter pairing neural network for person re-identification. In *CVPR*, 2014.
- [29] Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting batch normalization for practical domain adaptation. *arXiv preprint arXiv:1603.04779*, 2016.
- [30] Yu-Jhe Li, Fu-En Yang, Yen-Cheng Liu, Yu-Ying Yeh, Xiaofei Du, and Yu-Chiang Frank Wang. Adaptation and re-identification network: An unsupervised deep transfer learning approach to person re-identification. In *CVPR workshops*, 2018.
- [31] Shengcai Liao, Yang Hu, Xiangyu Zhu, and Stan Z Li. Person re-identification by local maximal occurrence representation and metric learning. In *CVPR*, 2015.
- [32] Shengcai Liao and Stan Z Li. Efficient psd constrained asymmetric metric learning for person re-identification. In *ICCV*, pages 3685–3693, 2015.
- [33] Liang Lin, Guangrun Wang, Wangmeng Zuo, Xiangchu Feng, and Lei Zhang. Cross-domain visual matching via generalized similarity measure and feature learning. *TPAMI*, 39(6):1089–1102, 2016.
- [34] Shan Lin, Haoliang Li, Chang-Tsun Li, and Alex Chichung Kot. Multi-task mid-level feature alignment network for unsupervised cross-dataset person re-identification. *BMVC*, 2018.
- [35] Jiawei Liu, Zheng-Jun Zha, Di Chen, Richang Hong, and Meng Wang. Adaptive transfer network for cross-domain person re-identification. In *CVPR*, 2019.
- [36] Chen Change Loy, Tao Xiang, and Shaogang Gong. Time-delayed correlation analysis for multi-camera activity understanding. *IJCV*, 90(1):106–129, 2010.
- [37] Hao Luo, Youzhi Gu, Xingyu Liao, Shenqi Lai, and Wei Jiang. Bag of tricks and a strong baseline for deep person re-identification. In *CVPR workshops*, 2019.

- [38] Ping Luo, Jiamin Ren, Zhanglin Peng, Ruimao Zhang, and Jingyu Li. Differentiable learning-to-normalize via switchable normalization. *ICLR*, 2019.
- [39] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *JMLR*, 2008.
- [40] Krikamol Muandet, David Balduzzi, and Bernhard Schölkopf. Domain generalization via invariant feature representation. In *ICML*, pages 10–18, 2013.
- [41] Xingang Pan, Ping Luo, Jianping Shi, and Xiaoou Tang. Two at once: Enhancing learning and generalization capacities via ibn-net. In *ECCV*, 2018.
- [42] Lei Qi, Lei Wang, Jing Huo, Luping Zhou, Yinghuan Shi, and Yang Gao. A novel unsupervised camera-aware domain adaptation framework for person re-identification. *ICCV*, 2019.
- [43] Xuelin Qian, Yanwei Fu, Wenxuan Wang, et al. Pose-normalized image generation for person re-identification. In *ECCV*, 2018.
- [44] Shiv Shankar, Vihari Piratla, Soumen Chakrabarti, Siddhartha Chaudhuri, Preethi Jyothi, and Sunita Sarawagi. Generalizing across domains via cross-gradient training. *ICLR*, 2018.
- [45] Jifei Song, Yongxin Yang, Yi-Zhe Song, Tao Xiang, and Timothy M Hospedales. Generalizable person re-identification by domain-invariant mapping network. In *CVPR*, 2019.
- [46] Liangchen Song, Cheng Wang, Lefei Zhang, Bo Du, Qian Zhang, Chang Huang, and Xinggang Wang. Unsupervised domain adaptive re-identification: Theory and practice. *arXiv preprint arXiv:1807.11334*, 2018.
- [47] Chi Su, Jianing Li, Shiliang Zhang, et al. Pose-driven deep convolutional model for person re-identification. In *ICCV*, 2017.
- [48] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *ECCV*, pages 480–496, 2018.
- [49] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.
- [50] Haotian Tang, Yiru Zhao, and Hongtao Lu. Unsupervised person re-identification with iterative self-supervised domain adaptation. In *CVPR workshops*, 2019.
- [51] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Instance normalization: The missing ingredient for fast stylization. *arXiv preprint arXiv:1607.08022*, 2016.
- [52] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. Improved texture networks: Maximizing quality and diversity in feed-forward stylization and texture synthesis. In *CVPR*, pages 6924–6932, 2017.
- [53] Jingya Wang, Xiatian Zhu, Shaogang Gong, and Wei Li. Transferable joint attribute-identity deep learning for unsupervised person re-identification. In *CVPR*, 2018.
- [54] Yan Wang, Lequn Wang, Yurong You, Xu Zou, Vincent Chen, Serena Li, Gao Huang, Bharath Hariharan, and Kilian Q Weinberger. Resource aware person re-identification across multiple resolutions. In *CVPR*, 2018.
- [55] Zhixiang Wang, Zheng Wang, Yinqiang Zheng, Yung-Yu Chuang, and Shin’ichi Satoh. Learning to reduce dual-level discrepancy for infrared-visible person re-identification. In *CVPR*, pages 618–626, 2019.
- [56] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer GAN to bridge domain gap for person re-identification. In *CVPR*, 2018.
- [57] Zheng Wei-Shi, Gong Shaogang, and Xiang Tao. Associating groups of people. In *BMVC*, pages 23–1, 2009.
- [58] Ancong Wu, Wei-Shi Zheng, Hong-Xing Yu, Shaogang Gong, and Jianhuang Lai. Rgb-infrared cross-modality person re-identification. In *ICCV*, pages 5380–5389, 2017.
- [59] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. End-to-end deep learning for person search. *arXiv preprint arXiv:1604.01850*, 2:2, 2016.
- [60] Qize Yang, Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Patch-based discriminative feature learning for unsupervised person re-identification. In *CVPR*, 2019.
- [61] Mang Ye, Xiangyuan Lan, Jiawei Li, and Pong C Yuen. Hierarchical discriminative learning for visible thermal person re-identification. In *AAAI*, 2018.
- [62] Mang Ye, Zheng Wang, Xiangyuan Lan, and Pong C Yuen. Visible thermal person re-identification via dual-constrained top-ranking. In *IJCAI*, pages 1092–1099, 2018.
- [63] Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Cross-view asymmetric metric learning for unsupervised person re-identification. In *ICCV*, 2017.
- [64] Hong-Xing Yu, Wei-Shi Zheng, Ancong Wu, Xiaowei Guo, Shaogang Gong, and Jian-Huang Lai. Unsupervised person re-identification by soft multilabel learning. In *CVPR*, 2019.
- [65] Li Zhang, Tao Xiang, and Shaogang Gong. Learning a discriminative null space for person re-identification. In *CVPR*, 2016.
- [66] Xinyu Zhang, Jiewei Cao, Chunhua Shen, and Mingyu You. Self-training with progressive augmentation for unsupervised cross-domain person re-identification. In *ICCV*, 2019.
- [67] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, et al. Densely semantically aligned person re-identification. In *CVPR*, 2019.
- [68] Haiyu Zhao, Maoqing Tian, Shuyang Sun, et al. Spindle net: Person re-identification with human body region guided feature decomposition and fusion. In *CVPR*, 2017.
- [69] Liang Zheng, Liyue Shen, et al. Scalable person re-identification: A benchmark. In *ICCV*, 2015.
- [70] Wei-Shi Zheng, Shaogang Gong, and Tao Xiang. Person re-identification by probabilistic relative distance comparison. In *CVPR*, 2011.
- [71] Zhedong Zheng, Liang Zheng, and Yi Yang. Unlabeled samples generated by gan improve the person re-identification baseline in vitro. In *ICCV*, 2017.
- [72] Zhun Zhong, Liang Zheng, Shaozi Li, and Yi Yang. Generalizing a person retrieval model hetero-and homogeneously. In *ECCV*, 2018.
- [73] Zhun Zhong, Liang Zheng, Zhiming Luo, Shaozi Li, and Yi Yang. Invariance matters: Exemplar memory for domain adaptive person re-identification. In *CVPR*, pages 598–607, 2019.

- [74] Zhun Zhong, Liang Zheng, Zhedong Zheng, Shaozi Li, and Yi Yang. Camstyle: A novel data augmentation method for person re-identification. *TIP*, 2018.
- [75] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, et al. Omni-scale feature learning for person re-identification. *ICCV*, 2019.

# Appendix

## 1. Implementation Details

**Network Details.** We use ResNet-50 [13, 1, 67, 37] as our base network for both baselines and our schemes. We build a strong baseline *Baseline* with some commonly used tricks integrated. Similar to [1, 67, 37], the last spatial down-sample operation in the last Conv block is removed. The proposed SNR module is added after the last layer of each convolutional block/stage of the first four stages. The input image resolution is  $256 \times 128$ .

**Data Augmentation.** We use the commonly used data augmentation strategies of random cropping [54, 67], horizontal flipping, and label smoothing regularization [49]. To enhance the generalization ability, we further incorporate some useful data augmentation tricks, such as color jittering and disabling random erasing (REA) [37, 75]. REA hurts models in cross-domain ReID task [37, 24], because REA which masks the regions of training images makes the model learn more knowledge in the training source domain. It causes the model to perform worse in the unseen target domain.

**Training Details for Domain Generalization.** Following [14], a batch is formed by first randomly sampling  $P$  identities. For each identity, we sample  $K$  images. Then the batch size is  $B = P \times K$ . We set  $P = 24$  and  $K = 4$  (*i.e.*, batch size  $B = P \times K = 96$ ).

We use the Adam optimizer [23] for model optimization. Similar to [37, 67], we first warm up the model for 20 epochs with a linear growth learning rate from  $8 \times 10^{-6}$  to  $8 \times 10^{-4}$ . Then we set the initial learning rate as  $8 \times 10^{-4}$  and optimize the Adam optimizer with a weight decay of  $5 \times 10^{-4}$ . The learning rate is decayed by a factor of 0.5 for every 40 epochs. Our model (here we use ResNet-50 as our backbone) with SNR converges well after training of 280 epochs and we use it for evaluating the generalization performance on target datasets. All our models are implemented on PyTorch and trained on a single 32G NVIDIA-V100 GPU.

**Training Details for Domain Adaptation.** For unsupervised domain adaptation person ReID, we combine our network with the unsupervised ReID approach MAR [64] for fine-tuning on the unlabelled target domain data. MAR [64] plays the role of assigning pseudo labels by hard negative mining, which facilitates the fine-tuning of base network. Similar to [64], during the fine-tuning, both source labeled data and target unlabelled data are jointly used for effective joint training. Specifically, during fine-tuning, a training

batch of size 96 is composed of 1) labeled source data (size  $B_1 = P \times K = 48$ , where  $P = 12, K = 4$ ), and 2) un-labeled target data (size  $B_2 = 48$ ). For the labeled source data, we optimize the network with the ReID loss  $\mathcal{L}_{ReID}$  and the proposed dual causality loss  $\mathcal{L}_{SNR}$ . For the un-labeled target data, we follow the adaptation strategy of MAR [64] to assign a pseudo soft multilabel for each sample and utilize these pseudo labels to perform soft multilabel-guided hard negative mining for training. We fine-tune the network also with the Adam optimizer [23] with a initial learning rate of  $1 \times 10^{-5}$  for 200 epochs. We optimize the Adam optimizer with a weight decay of  $5 \times 10^{-4}$ . The learning rate is decayed by a factor of 0.5 at 50, 100 and 150 epochs.

**Why do we perform disentanglement only on channel level?** We perform feature disentanglement only on channel level for **two reasons**: **1)** Those identity-irrelevant style factors (*e.g.*, illumination, contrast, saturation) are typically regarded as spatially consistent, which are hard to disentangle by spatial-attention. **2)** In our SNR, “disentanglement” aims at better “restitution” of the lost discriminative information due to Instance Normalization (IN). IN reduces style discrepancy of input features by performing normalization across spatial dimensions independently for each channel, where the normalization parameters are the same across different spatial positions. To be consistent with IN, we disentangle the features and reconstitute the identity-relevant ones to the normalized features on channel level.

## 2. Details of Datasets

Table 4: Details about the ReID datasets.

Datasets	Identities	Images	Cameras	Scene
Market1501 [69]	1501	32668	6	outdoor
DukeMTMC-reID [71]	1404	32948	8	outdoor
CUHK03 [28]	1467	28192	2	indoor
MSMT17 [56]	4101	126142	15	outdoor, indoor
VIPeR [11]	632	1264	2	outdoor
PRID2011 [15]	385	1134	2	outdoor
GRID [36]	250	500	2	indoor
i-LIDS [57]	119	476	N/A	indoor

In Table 4, we present the detailed information about the related person ReID datasets. Market1501 [69], DukeMTMC-reID [71], CUHK03 [28], and large-scale MSMT17 [56] are the most commonly used datasets for fully supervised ReID [67, 75] and unsupervised domain adaption ReID [64, 66, 8]. VIPeR [11], PRID2011 [15], GRID [36], and i-LIDS [57] are small ReID datasets which could be used for evaluating cross-domain/generalizable person ReID [45, 19, 24]. Market1501 [69] and DukeMTMC-reID [71] have pre-established test probe and test gallery splits which we use for our training and cross-test (*i.e.*,  $M \rightarrow D, D \rightarrow M$ ). For the smaller datasets (VIPeR, PRID2011, GRID, and i-LIDS), we use the standard 10 ran-

Table 5: Performance (%) comparisons of our scheme and others to demonstrate the effectiveness of our SNR module for generalizable person ReID. The rows denote source dataset(s) for training and the columns correspond to different target datasets for testing. We mask the results of supervised ReID by gray where the testing domain has been seen in training. Note that we show the total number of source training images by data num..

Source	Method	Target: Market1501		Target: Duke		Target: PRID		Target: GRID		Target: VIPeR		Target: iLIDs	
		mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1
Market1501 (M) data num. 32.6k	Baseline	82.8	93.2	19.8	35.3	13.7	6.0	25.8	16.0	37.6	28.5	61.5	53.3
	Baseline-A-IN	75.3	89.8	24.1	42.7	33.9	21.0	35.6	27.2	38.1	29.1	64.2	55.0
	Baseline-IBN	81.1	92.2	21.5	39.2	19.1	12.0	27.5	19.2	32.1	23.4	58.3	48.3
	Baseline-A-SN	83.2	93.9	20.1	38.0	35.4	25.0	29.0	22.0	32.2	23.4	53.4	43.3
	Baseline-IN	79.5	90.9	25.1	44.9	35.0	25.0	35.7	27.8	35.1	27.5	64.0	54.2
	<b>Baseline-SNR (Ours)</b>	<b>84.7</b>	<b>94.4</b>	<b>33.6</b>	<b>55.1</b>	<b>42.2</b>	<b>30.0</b>	<b>36.7</b>	<b>29.0</b>	<b>42.3</b>	<b>32.3</b>	<b>65.6</b>	<b>56.7</b>
Duke (D) data num. 32.9k	Baseline	21.8	48.3	71.2	83.4	15.7	11.0	14.5	8.8	37.0	26.9	68.3	58.3
	Baseline-A-IN	26.5	56.0	64.5	78.9	38.6	29.0	19.6	13.6	35.1	27.2	67.4	56.7
	Baseline-IBN	24.6	52.5	69.5	81.4	27.4	19.0	19.9	12.0	32.8	23.4	63.5	61.7
	Baseline-A-SN	25.3	55.0	73.0	85.9	41.4	32.0	18.8	12.8	31.3	24.1	64.8	63.3
	Baseline-IN	27.2	58.5	68.9	80.4	40.5	27.0	20.3	13.2	34.6	26.3	70.6	65
	<b>Baseline-SNR (Ours)</b>	<b>33.9</b>	<b>66.7</b>	<b>72.9</b>	<b>84.4</b>	<b>45.4</b>	<b>35.0</b>	<b>35.3</b>	<b>26.0</b>	<b>41.2</b>	<b>32.6</b>	<b>79.3</b>	<b>68.7</b>
Market1501 (M) + Duke (D) data num. 65.5k	Baseline	72.6	88.2	60.0	77.8	14.8	9.0	23.1	15.2	39.4	30.4	74.3	65.0
	Baseline-A-IN	76.5	91.4	62.2	80.1	45.0	30.0	36.7	28.0	37.3	28.2	73.6	65.2
	Baseline-IBN	74.6	90.4	62.3	80.1	43.7	32.0	32.6	24.0	42.8	33.2	73.8	65.0
	Baseline-A-SN	73.1	89.8	61.7	79.0	47.9	37.0	28.0	21.6	38.0	28.8	68.1	61.7
	Baseline-IN	77.5	91.6	63.9	81.5	48.1	36.0	39.2	31.2	43.8	33.9	73.2	64.3
	<b>Baseline-SNR (Ours)</b>	<b>80.3</b>	<b>92.9</b>	<b>67.2</b>	<b>83.1</b>	<b>57.9</b>	<b>50.0</b>	<b>41.3</b>	<b>34.4</b>	<b>46.7</b>	<b>37.7</b>	<b>85.2</b>	<b>80.0</b>
Market1501 (M) + Duke (D) + CUHK03 (C) data num. 93.7k	Baseline	76.4	89.8	63.6	79.0	27.0	19.0	25.7	18.4	46.3	36.4	77.1	66.3
	Baseline-A-IN	76.8	90.7	63.0	81.3	55.6	44.0	40.8	33.6	50.9	41.8	77.7	70.0
	Baseline-IBN	76.2	91.3	62.8	80.5	56.6	48.0	40.9	31.2	48.4	38.9	76.9	68.3
	Baseline-A-SN	71.1	89.3	62.0	78.8	55.4	46.0	34.1	26.4	50.3	39.8	79.6	71.7
	Baseline-IN	77.8	91.3	64.4	81.6	56.4	47.0	41.0	31.8	49.3	39.9	80.9	74.7
	<b>Baseline-SNR (Ours)</b>	<b>81.2</b>	<b>93.3</b>	<b>68.4</b>	<b>84.2</b>	<b>60.9</b>	<b>52.0</b>	<b>45.2</b>	<b>36.8</b>	<b>52.3</b>	<b>42.4</b>	<b>91.0</b>	<b>86.7</b>
MSMT17 (MT) data num. 126k	Baseline	23.1	48.2	29.2	47.6	16.4	11.0	9.8	5.6	40.8	30.1	74.0	66.7
<b>Baseline-SNR (Ours)</b>	<b>40.9</b>	<b>69.5</b>	<b>49.9</b>	<b>69.2</b>	<b>48.4</b>	<b>39.0</b>	<b>30.3</b>	<b>24.0</b>	<b>57.2</b>	<b>47.5</b>	<b>87.7</b>	<b>81.9</b>	
M + D + C + MT data num. 220k	Baseline	72.4	88.7	70.1	83.8	39.0	28.0	29.6	20.8	52.1	41.5	89.0	85.0
<b>Baseline-SNR (Ours)</b>	<b>82.3</b>	<b>93.4</b>	<b>73.2</b>	<b>85.5</b>	<b>60.0</b>	<b>49.0</b>	<b>41.3</b>	<b>30.4</b>	<b>65.0</b>	<b>55.1</b>	<b>91.9</b>	<b>87.0</b>	



Figure 5: Person images from different ReID datasets: Market-1501 [69], DukeMTMC-reID [71], CUHK03 [28], MSMT17 [56], and the four small-scale ReID datasets of PRID [15], GRID [36], VIPeR [11], and i-LIDS [57]. All images have been re-sized to 256×128 for easier comparison. We observe there are obvious domain gaps/style discrepancies across different datasets, especially for PRID [15] and GRID [36].

dom splits as in [19, 24] for testing (the four small datasets are not involved in training). CUHK03 [28] and MSMT17 [56] are used for training.

We randomly pick up 10 identities from each ReID

dataset and show them in Figure 5. We observe that: 1) there is style discrepancy across datasets, which is rather obvious for PRID and GRID; 2) MSMT17 has large style variants within the same dataset.

### 3. More Ablation Study Results

We show more comparisons of our scheme and others to demonstrate the effectiveness of our SNR module for generalizable person ReID in Table 5.

We have observations consistent with those in our paper. 1) IN-related baselines bring generalization ability improvement but decrease the performance for the same-domain. 2) Our *Baseline-SNR* achieves superior generalization capability thanks to the restitution of identity-relevant information by the SNR modules. 3) The generalization performance on unseen target domain increases consistently as the number of source datasets increases.

In Table 5, we also present the total number of source training images as marked by *data num.*  $N$ . For the single source dataset settings, MSMT17 is the largest dataset, which contains 126k images while Market1501 or Duke has about 33K images. For the target testing datasets VIPeR and iLIDs, the performance of *Baseline* trained by



Figure 6: Activation maps of our scheme (bottom) and the strong baseline *Baseline* (top) corresponding to images of varied styles. The maps of our method are more consistent/invariant to style variants.

this large scale dataset MSMT17 is 3.8% to 12.5% higher than those trained by Market1501 or Duke in mAP. Generally, the increase of training data could improve the performance. However, the performance of *Baseline* trained by MSMT17 has a rather low mAP accuracy of 9.8% on the target dataset GRID, being even poorer than that trained on Market1501 (25.8%) or Duke (14.5%). For the target dataset PRID, similarly, MSMT17 does not provide clear superiority. These indicate that it is not always true that a larger amount of training data results in better performance. The domain gap between MSMT17 and GRID is larger than that between Market1501/Duke and GRID. To validate this, we analyze the feature divergence (FD, detailed descriptions can be found in Section 4 below) between GRID and MSMT17, Market1501, Duke, respectively. We find that the divergence (here we calculate the feature divergence of the third convolutional block/stage within our *Baseline-SNR* trained by combining all the four datasets) of Market1501 vs. GRID, Duke vs. GRID, MSMT17 vs. GRID are 2.17, 3.49, and 4.51, respectively. Note that the larger the FD value, the larger the feature discrepancy between the two domains. The domain gap between MSMT17 and GRID is larger than that between Market1501 (or Duke) and GRID. For the similar reason, we find that additionally adding MSMT17 as the source training data does not bring further performance improvement on GRID and PRID target datasets in our scheme *Baseline-SNR* in comparison with the model trained by M+D+C source datasets.

#### 4. More Visualization Analysis

**More Feature Map Visualization.** In our paper, we compare the activation maps  $\tilde{F}^+$  of our scheme and those of the strong baseline scheme *Baseline* by varying the styles of input images (*e.g.*, contrast, illumination, saturation). Here, Figure 6(a) shows more visualization and Figure 6(b) shows visualization results on real images. We have the similar observations that the activation maps of our scheme are more consistent/invariant to style variants.

**Feature Divergence Analysis.** We analyze the feature di-

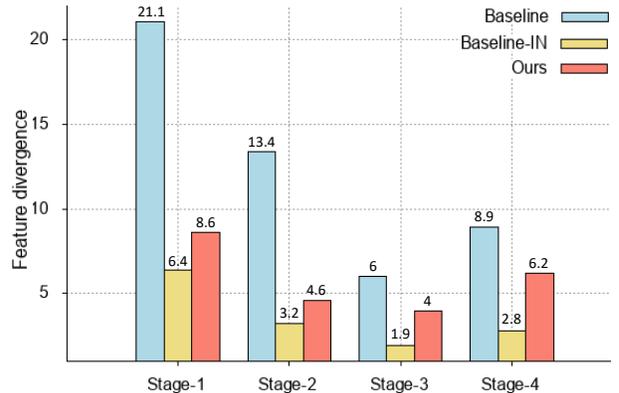


Figure 7: Analysis of the feature divergence between two different domains, Market1501 and Duke.

vergence between two datasets on three schemes: *Baseline*, *Baseline-IN*, and ours *SNR*, respectively. Following [41, 29], we use the symmetric KL divergence of features between domain A and B as the metric to measure feature divergence of the two domains. We train the models using Market1501 training dataset and evaluate the feature divergences between the test set of Market1501 and Duke (500 samples are randomly selected from each set). We calculate the feature divergence of the four convolutional blocks/stages respectively and show the results in Figure 7.

We observe that the feature divergence (FD) is large for *Baseline*. The introduction of IN as in scheme *Baseline-IN* significantly reduces the FD on all the four stages. The FD of Stage-4 is higher than that in Stage-3. That is likely because Stage-4 is more related to high-level discriminative semantic features for distinguishing different identities. The discrimination may increase the feature divergence. Due to the introduction of the SNR modules, the FD on all convolutional blocks/stages is also significantly reduced in our scheme in comparison with *Baseline*. It is higher than that of the scheme *Baseline-IN* which is probably because the restitution of some identity-relevant features increases the discrimination capability and thus increases the FD.

**Visualization of ReID Feature Vector Distributions.** In

Table 6: Performance (%) comparisons with the state-of-the-art approaches for the Domain Generalizable Person ReID (top rows) and Unsupervised Domain Adaptation for Person ReID (bottom rows), respectively. “(U)” denotes “unlabeled”. We mask the schemes of our *Baseline* and our *Baseline* with SNR modules (*i.e.*, *SNR(Ours)*) by gray, with fair comparison between each pair to validate the effectiveness of SNR modules.

Method	Venue	Source	Target: Duke		Source	Target: Market1501			
			mAP	Rank-1		mAP	Rank-1		
Domain Generalization (w/o using target data)	OSNet-IBN [75]	ICCV'19	Market1501	26.7	48.5	Duke	26.1	57.7	
	Baseline	This work	Market1501	19.8	35.3	Duke	21.8	48.3	
	Baseline-IBN [19]	BMVC'19	Market1501	21.5	39.2	Duke	24.6	52.5	
	<b>SNR(Ours)</b>	This work	Market1501	<b>33.6</b>	<b>55.1</b>	Duke	<b>33.9</b>	<b>66.7</b>	
	StrongBaseline [24]	ArXiv'19	MSMT17	43.3	64.5	MSMT17	36.6	64.8	
	OSNet-IBN [75]	ICCV'19	MSMT17	45.6	67.4	MSMT17	37.2	66.5	
	Baseline	This work	MSMT17	39.1	60.4	MSMT17	33.8	59.9	
	<b>SNR(Ours)</b>	This work	MSMT17	<b>50.0</b>	<b>69.2</b>	MSMT17	<b>41.4</b>	<b>70.1</b>	
	Unsupervised Domain Adaptation (using unlabeled target data)	PTGAN [56]	CVPR'18	Market1501 + Duke (U)	-	27.4	Duke + Market1501 (U)	-	38.6
		PUL [7]	TOMM'18	Market1501 + Duke (U)	16.4	30.0	Duke + Market1501 (U)	20.5	45.5
MMFA [34]		BMVC'18	Market1501 + Duke (U)	24.7	45.3	Duke + Market1501 (U)	27.4	56.7	
SPGAN [5]		CVPR'18	Market1501 + Duke (U)	26.2	46.4	Duke + Market1501 (U)	26.7	57.7	
TJ-AIDL [53]		CVPR'18	Market1501 + Duke (U)	23.0	44.3	Duke + Market1501 (U)	26.5	58.2	
ATNet [35]		CVPR'19	Market1501 + Duke (U)	24.9	45.1	Duke + Market1501 (U)	25.6	55.7	
CamStyle [74]		TIP'19	Market1501 + Duke (U)	25.1	48.4	Duke + Market1501 (U)	27.4	58.8	
HHL [72]		ECCV'18	Market1501 + Duke (U)	27.2	46.9	Duke + Market1501 (U)	31.4	62.2	
ARN [30]		CVPRW'19	Market1501 + Duke (U)	33.4	60.2	Duke + Market1501 (U)	39.4	70.3	
ECN [73]		CVPR'19	Market1501 + Duke (U)	40.4	63.3	Duke + Market1501 (U)	43.0	75.1	
UDAP [46]		ArXiv'18	Market1501 + Duke (U)	49.0	68.4	Duke + Market1501 (U)	53.7	75.8	
PAST [66]		ICCV'19	Market1501 + Duke (U)	54.3	72.4	Duke + Market1501 (U)	54.6	78.4	
SSG [8]		ICCV'19	Market1501 + Duke (U)	53.4	73.0	Duke + Market1501 (U)	58.3	80.0	
Baseline+MAR [64]		This work	Market1501 + Duke (U)	35.2	56.5	Duke + Market1501 (U)	37.2	62.4	
<b>SNR(Ours)+MAR [64]</b>		This work	Market1501 + Duke (U)	<b>58.1</b>	<b>76.3</b>	Duke + Market1501 (U)	<b>61.7</b>	<b>82.8</b>	
MAR [64]		CVPR'19	MSMT17 + Duke (U)	48.0	67.1	MSMT17 + Market1501 (U)	40.0	67.7	
PAUL [60]		CVPR'19	MSMT17 + Duke (U)	53.2	72.0	MSMT17 + Market1501 (U)	40.1	68.5	
Baseline+MAR [64]		This work	MSMT17 + Duke (U)	46.2	66.3	MSMT17 + Market1501 (U)	39.4	66.9	
<b>SNR(Ours) + MAR [64]</b>		This work	MSMT17 + Duke (U)	<b>61.6</b>	<b>78.2</b>	MSMT17 + Market1501 (U)	<b>65.9</b>	<b>85.5</b>	

Table 7: Performance (%) comparison with the latest domain generalizable ReID method Domain-Invariant Mapping Network (DIMN) [45] under the same experimental setting (*i.e.*, training on the same five datasets, Market1501[69]+DukeMTMC-reID[71]+CUHK02[27]+CUHK03[28]+CUHK-SYSU[59]).

Source	Method	Target: PRID		Target: GRID		Target: VIPeR		Target: iLIDs	
		mAP	Rank-1	mAP	Rank-1	mAP	Rank-1	mAP	Rank-1
Market + Duke + CUHK02 + CUHK03 + CUHK-SYSU	DIMN [45] CVPR'19	51.9	39.2	41.1	29.3	60.1	51.2	78.4	70.2
	Baseline	43.8	35.0	37.7	28.0	54.6	45.6	75.3	65.0
	<b>SNR (Ours)</b>	<b>66.5</b>	<b>52.1</b>	<b>47.7</b>	<b>40.2</b>	<b>61.3</b>	<b>52.9</b>	<b>89.9</b>	<b>84.1</b>

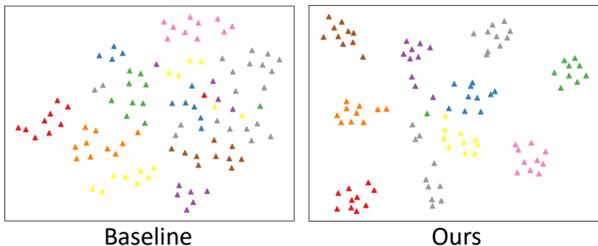


Figure 8: Visualization of the final ReID feature vector distribution for *Baseline* and *Ours* on the unseen target dataset Duke. Different identities are denoted by different colors.

Figure 8, we further visualize the distribution of the final ReID feature vectors using t-SNE [39] for *Baseline* scheme and our final scheme on the unseen target dataset Duke

Table 8: Differences between settings of supervised, domain adaptive, and domain generalizable ReID.

Setting	Use target domain data?	Use target domain label?
Supervised	✓	✓
Domain adaptation	✓	✗
Domain generalization	✗	✗

(*i.e.*, Market1501→Duke). In comparison with *Baseline*, the feature distribution of the same identity (same color) becomes more compact while those of the different identities are pushed away in our scheme. It is easier to distinguish between different identities by our method.

Table 9: Performance (%) comparisons with the state-of-the-art RGB-IR ReID approaches on SYSU-MM01 dataset. R1, R10, R20 denote Rank-1, Rank-10 and Rank-20 accuracy, respectively.

Method	Venue	All Search								Indoor-Search							
		Single-Shot				Multi-shot				Single-Shot				Multi-Shot			
		mAP	R1	R10	R20	mAP	R1	R10	R20	mAP	R1	R10	R20	mAP	R1	R10	R20
HOG [4]	CVPR'05	4.24	2.76	18.3	32.0	2.16	3.82	22.8	37.7	7.25	3.22	24.7	44.6	3.51	4.75	29.1	49.4
MLBP [32]	ICCV'15	3.86	2.12	16.2	28.3	-	-	-	-	-	-	-	-	-	-	-	-
LOMO [31]	CVPR'15	4.53	3.64	23.2	37.3	2.28	4.70	28.3	43.1	10.2	5.75	34.4	54.9	5.64	7.36	40.4	60.4
GSM [33]	TPAMI'17	8.00	5.29	33.7	53.0	-	-	-	-	-	-	-	-	-	-	-	-
One-stream [58]	ICCV'17	13.7	12.1	49.7	66.8	8.59	16.3	58.2	75.1	56.0	17.0	63.6	82.1	15.1	22.7	71.8	87.9
Two-stream [58]	ICCV'17	12.9	11.7	48.0	65.5	8.03	16.4	58.4	74.5	21.5	15.6	61.2	81.1	14.0	22.5	72.3	88.7
Zero-Padding [58]	ICCV'17	16.0	14.8	52.2	71.4	10.9	19.2	61.4	78.5	27.0	20.6	68.4	85.8	18.7	24.5	75.9	91.4
TONE [61]	AAAI'18	14.4	12.5	50.7	68.6	-	-	-	-	-	-	-	-	-	-	-	-
HCML [61]	AAAI'18	16.2	14.3	53.2	69.2	-	-	-	-	-	-	-	-	-	-	-	-
BCTR [62]	IJCAI'18	19.2	16.2	54.9	71.5	-	-	-	-	-	-	-	-	-	-	-	-
BDTR [62]	IJCAI'18	19.7	17.1	55.5	72.0	-	-	-	-	-	-	-	-	-	-	-	-
D-HSME [12]	AAAI'19	23.2	20.7	62.8	78.0	-	-	-	-	-	-	-	-	-	-	-	-
cmGAN [3]	IJCAI'18	27.8	27.0	67.5	80.6	22.3	31.5	72.7	85.0	42.2	31.7	77.2	89.2	32.8	37.0	80.9	92.3
D <sup>2</sup> RL [55]	CVPR'19	29.2	28.9	70.6	82.4	-	-	-	-	-	-	-	-	-	-	-	-
Baseline	This work	25.5	26.3	66.7	80.2	19.2	32.7	73.5	86.8	39.4	30.8	75.1	86.8	29.0	40.1	83.1	93.6
<b>Ours</b>	This work	<b>33.9</b>	<b>34.6</b>	<b>75.9</b>	<b>86.6</b>	<b>27.4</b>	<b>41.7</b>	<b>83.3</b>	<b>92.3</b>	<b>50.4</b>	<b>40.9</b>	<b>83.8</b>	<b>91.8</b>	<b>40.5</b>	<b>50.0</b>	<b>91.4</b>	<b>96.1</b>

## 5. Comparison with State-of-the-Arts (Complete version)

To save space, we only present the latest approaches in the paper and here we show comparisons with more approaches in Table 6. Besides the description in *Introduction* and *Related Work* sections of our paper, we illustrate the difference between domain generalization and domain adaptation for person ReID in Table 8.

Moreover, in Table 7, we further compare our *SNR* with the latest generalizable ReID method Domain-Invariant Mapping Network (*DIMN*) [45] under the same experimental setting, *i.e.*, training on the same five datasets, Market1501 [69] + DukeMTMC-reID [71] + CUHK02 [27] + CUHK03 [28] + CUHK-SYSU [59]. We observe that *SNR* not only outperforms the *Baseline* by a large margin (up to 22.7% in mAP on PRID), but also significantly outperforms *DIMN*[45] by **14.6%/6.6%/1.2%/11.5%** in mAP on PRID/GRID/ViPeR/i-LIDS, respectively.

## 6. Performance on Another Backbone

Our *SNR* is a plug-and-play module which can be added to available ReID networks. We integrate it into the recently proposed lightweight ReID network OSNet [75] and Table 10 shows the results. We can see that by simply inserting *SNR* modules between the OS-Blocks, the new scheme *OSNet-SNR* outperforms their best model *OSNet-IBN* by 5.0% and 5.5% in mAP for M→D and D→M, respectively. Note that, for fair comparison, we use the official released weights and codes <sup>2</sup> of OSNet [75] to conduct these experiments.

<sup>2</sup><https://github.com/KaiyangZhou/deep-person-reid>

Table 10: Evaluation of the generalization capability of proposed *SNR* modules on OSNet [75]. We use the official released weights and codes of OSNet for the experiments.

Method	M→D		D→M	
	mAP	Rank-1	mAP	Rank-1
Baseline (ResNet50)	19.8	35.3	21.8	48.3
OSNet [75]	19.3	35.2	21.7	49.9
OSNet-IBN [75]	26.7	48.5	26.1	57.7
<b>OSNet-SNR</b>	<b>31.7</b>	<b>53.6</b>	<b>31.6</b>	<b>62.7</b>

## 7. RGB-Infrared Cross-Modality Person ReID

To further demonstrate the generalization capability of the proposed *SNR* module, we conduct experiment on a more challenging RGB-Infrared cross-modality person ReID task, where there is a large style discrepancy between RGB images and Infrared images.

We evaluate our models on the standard benchmark dataset SYSU-MM01 [58]. Following [58], we conduct evaluation using the released official code based on the average of 10 repeated random split of gallery and probe sets. As shown in Table 9, in comparison with *Baseline*, our scheme which integrates the proposed *SNR* module on *Baseline* achieves a significant gain of **8.4%**, **8.2%**, **11.0%**, and **11.5%** in terms of mAP under 4 different experimental settings, and achieves the state-of-the-art performance.