

Pose Invariant Person Re-Identification using Robust Pose-transformation GAN

Arnab Karmakar, and Deepak Mishra, *Member IEEE*

Abstract—The objective of person re-identification (re-ID) is to retrieve a person’s images from an image gallery, given a single instance of the person of interest. Despite several advancements, learning discriminative identity-sensitive and viewpoint invariant features for robust Person Re-identification is a major challenge owing to large pose variation of humans. This paper proposes a re-ID pipeline which utilizes the image generation capability of Generative Adversarial Networks combined with pose clustering and feature fusion to achieve pose invariant feature learning. The objective is to model a given person under different viewpoints and large pose changes and extract the most discriminative features from all the appearances. The pose transformational GAN (pt-GAN) module is trained to generate a person’s image in any given pose. In order to identify the most significant poses for discriminative feature extraction, a Pose Clustering module is proposed. The given instance of the person is modelled in varying poses and these features are effectively combined through the Feature Fusion Network. The final re-ID model consisting of these 3 sub blocks, alleviates the pose dependence in person re-ID. Also, The proposed model is robust to occlusion, scale, rotation and illumination, providing a framework for viewpoint invariant feature learning. The proposed method outperforms the state-of-the-art GAN based models in 4 benchmark datasets. It also surpasses the state-of-the-art models that report higher re-ID accuracy in terms of improvement over baseline.

The code snippets¹ of this work can be found in <https://github.com/arnabk001/Pose-Invariant-Person-Re-Identification-using-Robust-Pose-Transformation-GAN>

Index Terms—Person Re-Identification, Pose Transformation, Generative Adversarial Networks, Pose Invariance.

I. INTRODUCTION

IN one of the earliest definitions [1], Person Re-Identification is described as, “to identify it (the person) as the same particular as one encountered on a previous occasion”. In image based re-ID, the detected and localised set of pedestrians construct the ‘gallery’ set. Given any query image, the re-ID model aims to retrieve all possible instances of the person-of-interest from the gallery set. Therefore, person re-ID is considered as an interdisciplinary field which sits in between image classification and information retrieval.

In recent years, person re-ID has gained significant importance due to its widespread applications in security and surveillance, public safety, crowd control etc. Combining

person re-ID with other branches of computer vision, such as face recognition, pose determination, activity recognition etc. head towards developing smart homes and cities. However, the primary bottleneck in developing robust person re-ID models arises from intra-class variations, i.e. the large appearance change of the same person across different scenes and camera views. This is mainly contributed to large pose change of the person-of-interest, which is also responsible for self-occlusion of the subject. The pose change of human subjects often lead to occlusion by external attributes (objects). In addition to that, in image-based re-ID, challenges such as inaccurate detection, occlusion, variation of illumination and resolution etc. are significant hindrances (Fig. 1) from achieving high accuracy with robustness.



Fig. 1. Challenges in ReID: (a-b) inaccurate detection, (c-d) pose misalignment, (e-f) occlusions, (g-h) very similar appearance of different IDs

Person re-identification is generally performed in two steps:

- 1) learning discriminative features from the given image of the person-of-interest, and
- 2) retrieve (and rank) all possible instances of the person-of-interest from the gallery set based on the feature similarity (distance metric)

Most existing methods [2] [3] make use of complex Deep Neural Networks (DNNs) as the backbone for feature extraction from pedestrian images. However, DNNs require a large number of samples (images) per person for proper training, hence it has limited success on smaller datasets. Also, the conventional problems of varying illumination, scale and occlusion are still persistent with DNN based model.

The idea of pose independence has been recently explored in the re-ID community [4] [5] [6] [7]. In this work, we argue that pose variation is the most significant hindrance when it comes to achieving highly accurate re-ID performance, and we focus on leveraging the pose dependence along with attaining invariance in occlusion, scale, illumination and rota-

A. Karmakar is with Human Space Flight Centre, Indian Space Research Organisation, Bengaluru, Karnataka 560231, India, e-mail: arnabkarmakar.001@gmail.com

D. Mishra is with Indian Institute of Space Science of Technology, Trivandrum, Kerala 695547, India, e-mail: deepak.mishra@iist.ac.in

Manuscript received March 19, 2020; revised April 26, 2021.

¹The code published here are snippets for various sub-parts of the whole project. Improvements are being made in the project and a fully executable code will be published soon.

tion. We achieve this by modelling the person-of-interest into multiple significant poses and then combining the extracted features from each of these models to generate a complete viewpoint invariant feature vector. The proposed end-to-end re-ID pipeline is constructed upon three building blocks:

- 1) A pose transformational GAN (pt-GAN) model to generate a given person image in any pose
- 2) An unsupervised Pose Clustering model to select the significant poses for pose transformation
- 3) A Feature Fusion Network (FusionNet) to effectively combine the discriminative features from the generated images and the source image.

At the heart of this re-ID model is the in-house developed pose transformational GAN (pt-GAN), which is an improvement over our previous work [8] in terms of loss function and training methodology; finetuned on re-ID datasets. Combining this with Pose Clustering and FusionNet, we emphasize our contribution in this paper as follows:

- 1) We propose an end-to-end feature extraction model for person re-ID while achieving pose invariance. Additionally, robustness is achieved in terms of illumination, scale, occlusion and rotation by using efficient data augmentation techniques and utilizing inherent CNN features through the GAN architecture.
- 2) The proposed end-to-end trainable model provides a powerful methodology to learn a viewpoint invariant feature representation from a single instance of the person-of-interest.
- 3) The proposed person re-ID framework incorporates an in-house developed pt-GAN architecture that serves as an efficient data augmentation tool. Incorporating this pose transformation GAN benefits in achieving an increase in the training dataset where the available dataset is small.

Our work is validated on four benchmark datasets, i.e. three large scale public datasets Market-1501 [9], CUHK03 [10] and DukeMTMC-reID [11] and a small-scale dataset CUHK01 [12]. **The proposed method surpasses the GAN based person re-ID models in both Rank-1 accuracy and mean Average Precision (mAP), and the results are comparable to the state-of-the-art methods. Additionally, the proposed method also outperforms the popular state-of-the-art re-ID models in terms of performance improvement over baseline.**

The rest of the paper is organised as follows: in section II, recent practices and developments in the re-ID community is explained, drawing a comparison to our method. Subsequently in section III, we formally define the problem of person re-identification. The intrinsic details of our model is demonstrated in section III and the experimental setting and detailed results are compared in section IV and V, respectively.

II. RELATED WORK

Initial approaches of re-ID started with handcrafted features [9] [13] [14], where a combination of manually extracted features are used as the image descriptor. Although these methods provide good accuracy in smaller datasets, there have

not been any significant improvement of these methods over the years. Another direction of work focus on developing improved distance metric learning methods [7] to construct discriminative subspace for the extracted features for better matching accuracy. However, the scope of distance metric in improving re-ID performance is limited when the dataset is large, unless a powerful feature extractor is employed. A comprehensive overview of the history and development of re-ID methodologies is described by Zheng et al. [15].

With the release of larger datasets and higher computational capability, Deep Learning based methods have provided significant performance improvement in recent years [15] [16]. A large portion of the existing methods depend heavily on the CNN classification backbone such as ResNet [17], DenseNet etc. In this transfer learning setting, the backbone CNN is finetuned using the re-ID data using the traditional classification loss [2], or using triplets in a Siamese learning setting [18] [19]. Some of the methods experiment with a part-based feature extraction and matching in addition with full image features [18]. Any of these methods, however, do not solve the problems of occlusion, background clutter, inaccurate detection, varying scale etc. Since deep learning methods give a holistic feature including the background, many methods use foreground segmentation or body-part based semantic segmentation to improve performance [5], which leads to computational overhead. Many attempts of combining the handcrafted features with the CNN based features have also been observed [4] [20] since deep learning based method are not very successful in smaller datasets.

A. Pose based methods in re-ID

Exploring the human body pose to improve re-ID performance has been a recent development. Many previous works have analysed the dependence of human pose in person re-ID [4] [5] [6] [7] In earlier works, the Symmetry-Driven Accumulation of Local Features (SDALF) method was proposed by Farenzena et al. [13] where handcrafted complementary features are extracted depending on localization of perceptually relevant human parts. They achieve pose invariance by combining these features through a weighting scheme based on human body (vertical) axis of symmetry. Weinrich et al. [14] detect the upper body pose for tracking the human subject in egocentric image frames. The local texture features of the upper body are learned and a generative 3D shape model is used for re-identification purpose. In another direction of research, there have been attempts to apply metric learning, based on a transformation function for different camera pairs, that can be used to model the pose change of the human and then compute feature distances to boost person re-ID. Following this philosophy, Bak et al. [7] try to learn a metric pool, i.e. a pose change metric based on Mahalanobis distance by classifying human pose into 3 groups: front, back and side. All the aforementioned methods use typical handcrafted features to achieve pose variation. However, in recent times there have been increasing application of deep learning based feature extractors and pose detectors for re-ID. Zheng et al. [7] addresses the pose misalignment problem by ‘PoseBox’, a

pose based body part extraction and reconstruction algorithm for the pedestrian images. The original image and the PoseBox output both are trained in a CNN to extract pose invariant features for re-ID. Cho et al. [21] proposes a multi-shot re-ID framework by learning the transformation matrix from one pose to another by utilizing external camera parameters to predict the next pose (target pose) of the pedestrian and then use the matching score for re-ID framework. Sarfraz et al. [6] uses the detected pose keypoints with the pedestrian image as the combined input to the CNN for learning pose sensitive embedding. The authors use a view predictor, i.e. pose classification branch in conjunction with the CNN branch for robust embedding generation. However, the improvement of the framework is largely contributed to the new re-ranking methodology. Zhao et al. [22] proposes the ‘SpindleNet’ architecture that uses a Region Proposal Network (RPN) utilizing the human body structure (pose), and then extract features from each body part using the ROI pooling. These part-based features are merged using a tree-like ‘competitive’ feature fusion network with the full-image feature. Following this idea, Jhonson et al. [4] tackles the person detection error by a SpindleNet-like body part based deep feature learning methodology by using an ensemble of handcrafted features such as Hue-Saturation-Value (HSV), Scale Invariant Local Ternary Pattern (SILTP) with the deeply learned model to produce better results. The fusion of Deep learning based features with handcrafted features have also been observed in Lee et al.’s [20] work, where the authors use an ensemble of invariant features for re-ID, by combining holistic (deep) features and regional (handcrafted) features. However, the authors do not provide specific justification on how these features are pose invariant.

B. GANs in Person re-ID

Goodfellow et al. [23] introduced GANs for image generation from random noise. Recently, the idea of utilizing GANs for improving person re-ID has been explored by the re-ID community. In the initial methods, GANs were used as a data augmentation technique to provide sufficient training data to strengthen the learning process. Zheng et al. [11] achieve an improvement over baseline using the GAN generated data and label smoothing regularization to train the CNN model in a semi-supervised setting. Liu et al. [24] proposes a pose transferrable GAN using a Generator-Guider-Discriminator architecture. A U-net [25] based structure was adopted for generator and the vanilla CNN structure was followed for the guider-discriminator model. The authors randomly sample poses from the MARS dataset and generate samples with new poses to increase the dataset size to finetune re-ID performance using the backbone CNN network.

Two significant works that incorporate pose invariance using GANs in person re-ID, are Qian et al. [26] and Ge et al. [27]. Ge et al. [27] propose a Siamese learning structure using identity discriminator and pose discriminator as well as verification loss. This helps in learning identity related and pose unrelated feature embedding which they utilise through the backbone (ResNet-50 [17]) network during feature generation

while testing. Qian et al. [26] train a pose transformation GAN and for every query image, generate 8 canonical pose transformed images. It uses ResNet [17] architecture to extract image features from the source image and generated images. However, the authors use 8 canonical poses for every image and combine the feature vectors by element wise maximum operation. Both these decisions are fairly crude and the authors do not provide enough justification for using this type of a learning setting. In our work, we explore the effect of the number of generated images (with new poses), and apply a fully connected neural network FusionNet for effective combination of features. Our results on several benchmark datasets validate this claim and we achieve the best results among all GAN based models for re-ID.

III. METHODOLOGY

A. Definition

Considering a closed world model, where a total of n images define the gallery set $\mathbb{G} : \{g_j\}_{j=1}^n$ having \mathcal{N} different identities. Given a query (probe) image q_i with identity i , the person reidentification framework aims to retrieve all possible instances of the same identity i from the given gallery set \mathbb{G}

$$\mathcal{R}_i^* = \arg \max_{i \in 1, 2, \dots, \mathcal{N}} \text{sim}(q_i, \{g_j\}_{j=1}^n) \quad (1)$$

where \mathcal{R}_i^* is the retrieved set of all images with the same identity i denoted by the probe image q_i and $\text{sim}(\cdot, \cdot)$ denotes some form of similarity function.

Assuming a training dataset of N images, denoted by $\mathcal{D}_{tr} = \{I_i, y_i\}_{i=1}^N$, where I_i denotes the person image and y_i denotes its identity, the idea is to learn a feature extraction function Φ in order to represent a given image I_i as a feature vector $\mathbf{f}_{I_i} = \Phi(I_i)$.

During testing, we need to retrieve all possible images of the same identity as the given probe image I_q from the testing (gallery) set $\mathcal{D}_{te} = \{I_j\}_{j=1}^{N'}$ based on the feature similarity measure $\text{sim}(\mathbf{f}_{I_q}, \{\mathbf{f}_{I_j}\}_{j=1}^{N'})$.

B. Framework Overview

During training, the input image I_i is transformed using ResNet50-1 model $\mathbf{F}_{R1}(\cdot)$ for the feature vector representation $\mathbf{F}_{R1}(I_i)$. This ResNet50-1 model is essentially a ResNet50 model trained on the ImageNet dataset.

The pose of I_i is detected by an off-the-shelf pose detection model Openpose [28] and the pose vector is denoted as p_i . The pose clustering module analyses all detected poses in the dataset to generate a set of K significant body poses

$$\mathcal{P}_K = \{p'_k\}_{k=1}^K = \{p'_1, p'_2, \dots, p'_K\}$$

The term p'_k denotes the k^{th} generated pose. The generator takes the input image feature $\mathbf{F}_{R1}(I_i)$ and the set of poses \mathcal{P}_K to generate K images,

$$\begin{aligned} \hat{I}_{i, \mathcal{P}_K} &= \mathbf{G}(\mathbf{F}_{R1}(I_i), \mathcal{P}_K) \\ \implies \hat{I}_{i, p'_k} \Big|_{k=1}^K &= \mathbf{G}(\mathbf{F}_{R1}(I_i), p'_k) \Big|_{k=1}^K \end{aligned}$$

$\hat{I}_{i, \mathcal{P}_K}$ denotes the original image I_i is being transformed to K poses. Now, the ResNet50-2 model $\mathbf{F}_{R2}(\cdot)$ is used to transform the K generated images into feature vectors $\mathbf{F}_{R2}(\hat{I}_{i, \mathcal{P}_K})$.

The ResNet50-1 model is pre-trained on the ImageNet dataset, whereas the ResNet50-2 model is finetuned on the re-ID dataset. Initially while training the pt-GAN model, the general image features of the person is of primary concern, and preserving the image features as a whole would give a better reconstruction. Hence a vanilla ResNet50-1 is deployed to extract the general features of the person. In the second stage, since the features of the original image along with generated images is more important for re-ID matching, we have used the ResNet50-2 model which is finetuned for re-ID purpose.

The K image feature vectors along with the input image feature $\mathbf{F}_{R1}(I_i)$ is transformed using the FusionNet model $\mathbf{F}_{FN}(\cdot)$ as the final feature vector. The complete transformation model, in accordance with the previously described notation, becomes

$$\mathbf{f}_{I_i} = \Phi(I_i) = \mathbf{F}_{FN}(\mathbf{F}_{R1}(I_i), \mathbf{F}_{R2}(\hat{I}_{i, \mathcal{P}_K})) \quad (2)$$

This feature vector \mathbf{f}_{I_i} extracted from all images in the gallery set, construct the feature space, where image retrieval is performed based on a similarity function $sim(\cdot, \cdot)$. The complete architecture is depicted in Figure 2.

C. Pose Guided Person Image Synthesis

In general, Generative Adversarial Networks (GANs) comprise of two parts: the Generator and the Discriminator. The Generator learns to approximate the underlying real data distribution to a random distribution, where sampling is performed to generate new samples. The Discriminator learns to distinguish the real life data and the generated data, thereby providing a measure of ‘goodness’ on how well the generator can mimic the real data. They both learn in a competitive manner.

In the pt-GAN model, we want the generator to learn the pedestrian image distribution, conditioned on human body pose. Our pt-GAN model is given an image pair $\{I_i, I_j\}$ of the same identity but different pose. The OpenPose [28] pose detection framework is used off-the-shelf to generate the pose vector p_j of the target image I_j . The generator takes in the image feature vector $\mathbf{F}_{R1}(I_i)$, concatenated with the target pose vector p_j to model the person in the desired pose $\hat{I}_j = \mathbf{G}(\mathbf{F}_{R1}(I_i), p_j)$

The generator learns the data distribution to minimize the difference between the generated image \hat{I}_j and the target image I_j . The loss function of the GAN model is denoted by,

$$\mathcal{L}(\mathbf{D}, \mathbf{G}) = \mathbb{E}_{I_j \sim p_{data}(I_j)} [\log(\mathbf{D}(I_j))] + \log(1 - \mathbf{D}(\mathbf{G}(\mathbf{F}_{R1}(I_i), p_j))) \quad (3)$$

We also use the added L2 norm for better reconstruction and cleaner image generation to the generator loss function.

$$\mathcal{L}_{L2} = \mathbb{E}_{I_j \sim p_{data}(I_j)} \|I_j - \hat{I}_j\|_2 \quad (4)$$

For Discriminator, we have added an extra classification branch in order to predict the class labels in addition to the

real/fake classification. Incorporating classification loss in the discriminator complements the generator’s capability to reconstruct detailed images with lesser artifacts. The categorical crossentropy loss is used here,

$$\mathcal{L}_c = - \sum_{i=1}^n \sum_{c=1}^C y_c \log(y'_c) \quad (5)$$

where y'_c is the predicted probability of class c , y_c is the actual class label; summed over all classes across all samples. Therefore we write the final loss function of our pt-GAN model as

$$\mathcal{L}(\mathbf{D}, \mathbf{G}) = \mathbb{E}_{I_j \sim p_{data}(I_j)} [\log(\mathbf{D}(I_j))] + \mathcal{L}_c + \log(1 - \mathbf{D}(\mathbf{G}(\mathbf{F}_{R1}(I_i), p_j))) + \mathcal{L}_{L2} \quad (6)$$

The GAN model is depicted in Figure 3. This is an improvement over our previous work [8] in terms of loss function formulation and training methodology, tuned for better performance in re-ID datasets. This model have been proven to produce robust human images with significant details. The image generation results are shown in section V-A.

D. Pose Clustering

After our pt-GAN model is trained to generate pedestrian images in any given pose, the requirement is to select the minimum number of best possible poses to synthesize the least number of images but still be able to extract viewpoint invariant features. Clustering techniques are applied to study the maximally occurring poses in the dataset, as they provide an unsupervised approach to structure the data points by extracting meaningful insights about the distribution. In this work, we have implemented clustering on the pose vectors based on two criteria, (1) Full-Body Pose based clustering, (2) Body joint based clustering method. For each of these methods, we implement two major clustering algorithms: (i) K-means Clustering and (ii) Gaussian Mixture Model based Clustering (GMM Clustering). The K-means method produces discretized output points while GMM produces a distribution of data points, where sampling is performed to produce output poses.

1) *Full-body pose based Clustering*: The complete pose vector p_i is taken as a single data point and then clustered into n_{cp} clusters. The K-means method directly produces the desired number of poses as the cluster centers. For image generation, these cluster centers are taken as the sample pose. Meanwhile the GMM method provides a pose distribution, and the required number of poses are randomly sampled and fed to the subsequent pt-GAN model to synthesize the pedestrian images.

2) *Body joint based Clustering*: In this approach, each body joint location is taken individually, and clustering is performed in every body joint position with n_{cbj} clusters, for j^{th} body joint position. For K-means, the cluster centers are selected randomly out of n_{cbj} clusters, for each body-joint, and the output pose is constructed as the accumulation of these body-joint co-ordinates. For GMM, full body pose is constructed using individual body-joint co-ordinate that are randomly sampled from the distribution of individual body-joints.

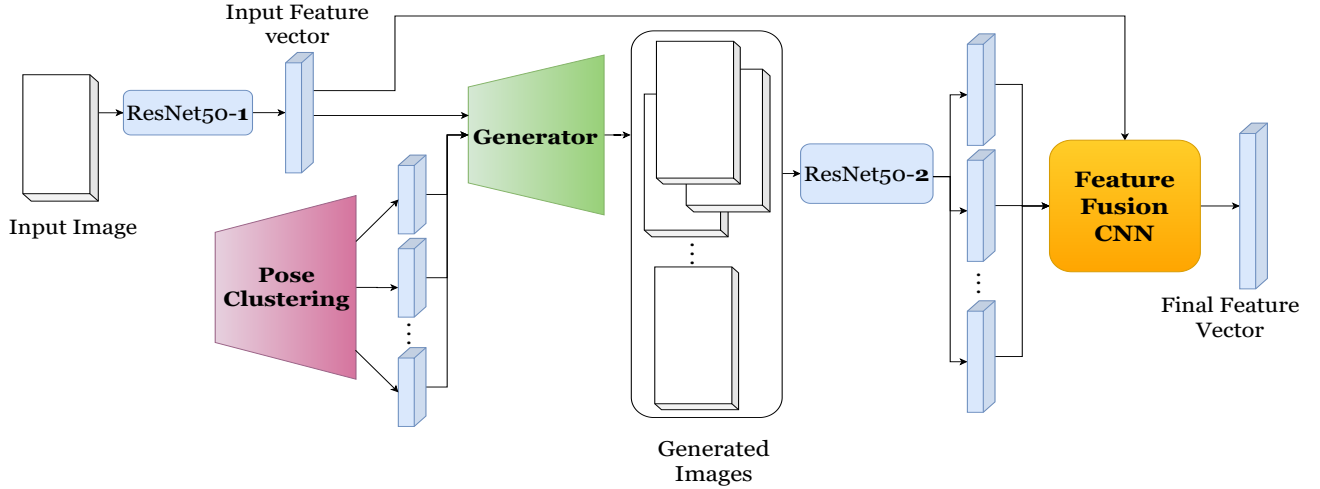


Fig. 2. Full architecture of the proposed person reID pipeline, consisting of the trained generator from our pt-GAN model, pose clustering and FusionNet. The number of generated images vary according to the pose clustering algorithm. The FusionNet is trained keeping the Generator and the Pose Clustering module fixed. During testing the features are extracted from second last layer of FusionNet.

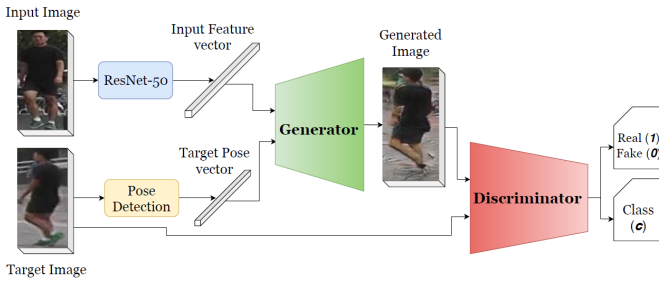


Fig. 3. The complete pt-GAN model architecture. The Generator learns pedestrian features from the image descriptors extracted by the ResNet50-1 model and reconstructs the pedestrian into the target pose vector. An extra classification branch is added to the Discriminator that complements the Generator’s capability to reconstruct detailed images with lesser artifacts.

E. Feature Fusion Network

The Pose Clustering module gives the N_{aug} best possible poses for the given image, and the pt-GAN module transforms the input image into these poses. Having these many samples augmented from one image, we need an efficient way of combining the best possible features of every image. We propose a Feature Fusion Network (FusionNet) to combine the features of the input image with all the augmented image.

As depicted in the architecture of the FusionNet in Figure 4, the image descriptors of the input image as well as the generated images are obtained using a ResNet50-2 model trained (finetuned) on re-ID dataset. We first perform concatenation of all the image descriptors to obtain an $(N_{aug} + 1)$ length input feature vector which is the input to the fully connected FusionNet. Since the input image holds the most prominent features for re-identification and the generated images might incorporate various distortions, we introduce a skip connection with the input feature to re-use the input features and imply more importance to the original image, as shown in Figure 4. During training, the final layer of this network is modified to predict the class labels and this model is trained on classification loss (categorical crossentropy).

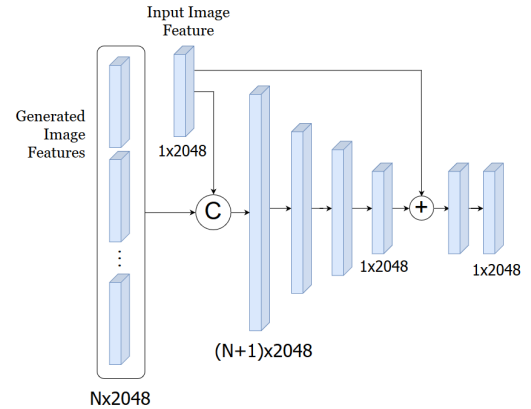


Fig. 4. Architecture of the proposed FusionNet. N stands for the number of generated images (correspondingly image feature vectors), $N \in \{8, 12, 16, 24\}$. The last layer of FusionNet is modified to the number of classes during training, while the features are extracted from second last layer during testing.

IV. EXPERIMENTS

A. Datasets

1) *Market-1501* [9]: This dataset contains 32,668 images of 1,501 identities, captured using 6 different cameras. We follow the standard split [9], where 12,936 images of 751 identities are used for training and the rest 19,732 images of 750 identities are used for testing. This dataset (training set) is used to train our pt-GAN model and the FusionNet model.

2) *DukeMTMC-reid* [11]: This dataset contains 36,411 images of 1,812 identities. We use 702 identities as the training set and the remaining 1,110 identities as the testing set, following the standard evaluation protocol [11].

3) *CUHK03* [10]: This dataset is captured using 6 cameras, containing 14,096 images of 1,467 identities. For training, 1367 identities are used. Validation and testing set consists of 100 identities each. We follow [10] for the testing process with 20 random splits.

4) *CUHK01 [12]*: This is a small scale dataset compared to the other 3, having 971 identities with 2 images per person. As in [12], the images of camera A is used as the probe image and camera B images are used as the gallery set.

B. Evaluation Metrics

The Cumulative Matching Curve (CMC) is generally used for evaluating the re-ID accuracy. However, for most of the existing methods the full details of the CMC curve is not shown, rather the rank-1, rank-5 and rank-10 accuracies are reported. Zheng et al. [9] proposes to use the mean Average Precision (mAP) as a single measure for overall correctness of re-ID models. In this work, we report the r-1, r-5 and r-10 accuracies for all 4 datasets, and the mAP for Market-1501 [9], DukeMTMC-reID [11] and CUHK03 [10].

C. Implementation details

1) *Data Augmentation*: Following [8], the image augmentation techniques are used as is. Specifically, Random Erasing [29] is used for occlusion invariance, and Random Crop is used as an augmentation for inaccurate detections. Also, random rotation and crop in the range $[+20 \text{ deg}, -20 \text{ deg}]$ is applied during the training of pt-GAN. Other methods such as random color jitter in all 3 channels, horizontal flip and random distortion are also incorporated while training the GAN model.

2) *GAN training*: For stable training and improving the performance of GANs, label smoothing is used with uniform random noise for the discriminator. The LeakyRelu activation and Adam optimizer (with $\beta_1 = 0.5$ and $\beta_2 = 0.999$) is used for both Generator and Discriminator. We have used the Kaiming-Normal initialization [17] for the GAN model and FusionNet. The other parameters such as the number of Residual blocks, the learning rate and the batch size is kept the same as our previous work [8].

The re-ID images were resized to 128×64 for training. The number of poses \mathcal{P}_K generated using the pose clustering module is varied from 8 to 24, i.e. $K \in \{8, 12, 16, 24\}$, and the ablation studies are presented in Table II. For full body pose based clustering, the number of clusters is equal to the number of generated poses ($n_{cp} = K$), whereas for body joint based clustering, each body joint is clustered into 3 clusters ($n_{cbj} = 3$, for each bodyjoint $j \in \{0, 1, 2, \dots, 24\}$). For K-means, one out of 3 coordinates are randomly selected while generating full body pose, where for GMM, it is sampled from the distribution.

3) *Model Parameters*: The feature size is taken as 2048 (output of the global average pooling layer of ResNet50) and the same is also replicated in the FusionNet architecture where the concatenated features are taken as input. The details of the FusionNet architecture is provided in Table I.

The FusionNet model contains a large number of trainable parameters. In order to prohibit overfitting, we have used dropout= 0.6 along with Weight Regularization, Batch Normalization and Early Stopping. Both the pt-GAN model as well as the pose clustering module is trained independently. After training, the FusionNet model is trained for classification accuracy keeping the other two models frozen.

TABLE I
DETAILS OF FUSIONNET ARCHITECTURE. FOR THE OPTIMAL CASE
GMM-12, $N = 12$.

Layer Name	Input feature size	Connected to	Output feature size	Params (M)
input_img	-	ResNet50-1	2048	-
gen_img_ft	-	ResNet50-2	$N*2048$	-
concat_1	2048, $N*2048$	input_img, gen_img_ft	$(N+1)*2048$	-
fc_1	$(N+1)*2048$	concat_1	$4*2048$	$(N+1)*16.7$
fc_2	$4*2048$	fc_1	2048	16.7
add_1	2048	fc_2, input_img	2048	-
output	2048	add_1	2048	4.1

V. RESULTS AND DISCUSSION

A. GAN Image Generation

The proposed GAN model is trained on re-ID datasets to generate a pedestrian image into any given pose, given only a single instance of the person-of-interest. Figure 5 demonstrates the generated images from our pt-GAN model. The original image of the pedestrian is shown in the left, and the generated images with the target poses are shown in the right. As it can be seen, our pt-GAN model is able to correctly identify and extract the significant features from the input image and then reconstruct the pedestrian appearance in the target pose. Features such as clothing, hair, male/female attributes are properly retained.

B. Pose Clustering

The pose clustering model provides a variety of poses to augment the input image. Figure 6 shows the poses obtained using K-means clustering in the Market-1501 dataset. Although K-means produces a fixed number of poses (according to the specified number of cluster centers), we achieve randomness when GMM is applied. Figure 7 shows a sample of generated images using GMM. In the final model, we have used GMM as it is able to provide superior results with less number of poses.

C. Re-ID performance

1) *Ablation studies*: We compare the results of both pose based and body-joint based clustering algorithms with number of poses varying from 8 to 24. As seen in Table II, K-means provides consistent performance in both the methods. The performance of GMM drops in body-joint based clustering algorithm as the distribution of each body joint is significantly sparse, when computed across all samples. We have observed the highest rank-1 accuracy in K-means-16, but the GMM-12 has been selected as the final model as it achieves higher mAP and comparable rank-1 accuracy with less number of poses.

2) *Results on benchmark Datasets*: Our results are compared against the state-of-the-art re-ID models. The Market-1501 and DukeMTMC-reID has been widely used in the literature to validate re-ID performance, but most works only report the rank-1 accuracy and mAP score, and the same approach



Fig. 5. Generated Image of our pt-GAN model. The original image is shown in the left. The generated images along with the given target pose is shown subsequently.

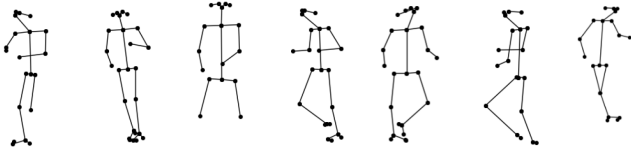


Fig. 6. Sample Poses obtained using K-means Clustering with full-body pose based method.

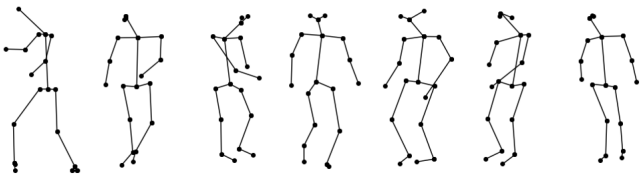


Fig. 7. Sampled Poses obtained using GMM with full-body pose based method

TABLE II
ABLATION STUDY OF THE PROPOSED MODEL ON MARKET-1501 DATASET. THE GMM-12 HAS BEEN SELECTED AS THE FINAL MODEL AS IT ACHIEVES HIGHER MAP AND VERY SIMILAR RANK-1 ACCURACY WITH LESS NUMBER OF POSES COMPARED TO K-MEANS-16.

Method		# poses	r-1	mAP
Pose based	Kmeans	8	86.40	80.23
		12	88.46	81.59
		16	88.40	82.10
		24	88.10	82.11
	GMM	8	88.21	81.61
		12	90.87	82.67
		16	90.81	82.43
		24	89.98	82.44
Body-joint based	Kmeans	8	86.74	81.20
		12	89.97	82.14
		16	90.88	81.45
		24	90.86	82.04
	GMM	8	71.41	68.42
		12	73.49	70.14
		16	74.13	70.89
		24	74.89	71.43

is being followed for the sake of comparison. However, for CUHK03 dataset, elaborate results containing rank-1, rank-5 and rank-10 accuracy alongwith mAP score is compared. Similar approach is taken while comparing our results in CUHK01 dataset where we compare our results on rank-1, rank-5 and rank-10 accuracy.

The accuracy comparison in Table III, IV, V and VI establish that our model is on par with the state-of-the-art methods. Our model has shown significant success not only in large datasets, but in smaller datasets (CUHK01) too. Also, it consistently outperforms all the GAN based methods developed for re-ID [26] [27] [11] [24]. The significant contribution of this work is the increase over the baseline. We have achieved an improvement of $\sim 9.64\%$ (rank-1) and $\sim 16.03\%$ (mAP) in the Market-1501 dataset, which is significantly higher compared to the existing works. In view of only rank-1 accuracy, the gain over baseline is $\sim 12.80\%$, $\sim 10.73\%$ and $\sim 13.96\%$ in DukeMTMC-reID, CUHK03 and CUHK01 dataset, respectively.

Table VII shows the relative improvement over baseline for the state-of-the-art models that have reported higher rank-1 accuracy values than the proposed model. The recent works show a higher accuracy mostly due to the incorporation of external features such as attention maps, semantic parsing etc. Most of these works follow a different philosophy and use a higher baseline while developing the model, thereby achieving higher accuracy. However, the proposed model excels in terms of percentage increment over baseline, primarily due to two factors:

- 1) Data augmentation for the pt-GAN model to gain independence from rotation, occlusion, illumination and scale; and
- 2) Feature extraction by combining the multi-pose (i.e. multi-view) images, generated by the pt-GAN model.

The improvement is also reflected in the retrieved images.

Figure 8 shows the retrieved images for 2 sample query image for ResNet50-1, ResNet50-2 and the proposed method. Since the ResNet50-2 model has been finetuned on reid datasets, its performance is better compared to vanilla ResNet50-1. However, the Proposed model shows the best accuracy among the three.

TABLE III

COMPARISON OF RESULTS IN THE MARKET-1501 DATASET. THE BEST AND SECOND BEST RESULTS ARE DENOTED WITH RED AND BLUE, RESPECTIVELY.

Method	Market-1501			
	Single Query		Multi Query	
	r-1	mAP	r-1	mAP
BOW [9]	44.42	20.76	-	-
LSTM Siamese [30]	-	-	61.60	35.31
Gated Siamese [19]	65.88	39.55	76.50	48.50
SpindleNet [22]	76.90	-	-	-
HP-Net [31]	76.90	-	-	-
PIE (poseBox) [5]	79.33	55.95	-	-
DLPAR [32]	81.00	63.40	-	-
SSM [33]	82.21	68.80	88.2	76.2
PDC [34]	84.14	63.41	-	-
SVDNet (RE) [29]	87.08	71.31	-	-
DML [35]	87.70	68.80	-	-
DeepTransfer [36]	83.70	65.50	89.60	73.80
JLML [37]	85.10	65.50	89.70	73.80
MLFN [38]	90.00	74.30	-	-
HA-CNN [39]	91.20	75.70	-	-
PCB [14]	93.80	81.60	-	-
US-GAN [11]	83.97	66.07	88.42	76.10
Liu et al. [24]	87.65	68.92	-	-
PN-GAN [26]	89.43	72.58	92.93	80.19
FD-GAN [27]	90.50	77.70	-	-
ResNet50-1 (Baseline-1)	83.69	76.48	-	-
ResNet50-2 (Baseline-2)	87.24	80.19	-	-
Ours	90.87	82.67	92.98	88.32
Ours (rerank)	91.76	88.74	93.64	89.24

3) *Robustness*: Our pt-GAN model shows good robustness towards occlusion, illumination and scale. As seen in the person 1 of Figure 5, the occlusion (railings) present in the input image is not propagated in any of the generated images. Fig. 9 demonstrates the occlusion invariance in the proposed model. The occluded images are given as the input to the model and are successfully reconstructed in the desired poses. Also, our model can handle human detection errors or human crop errors i.e. scaling issues. As depicted in person 2 of Figure 5, the pt-GAN model is able to correctly locate the human in the frame to extract only person-specific features for reconstruction.

A comparative study of feature distances is carried out to understand the intra-class and inter-class variations of extracted features in Figure 10. The intra-class feature distances is significantly lower, even in the presence of rotation and illumination variation as well as pose change; compared to that of the inter-class feature distances. Therefore, we achieve invariance in terms of pose, occlusion, scaling, rotation and illumination with our end-to-end re-ID model.

TABLE IV

COMPARISON OF RESULTS IN THE DUKEMTMC-reID DATASET. THE BEST AND SECOND BEST RESULTS ARE DENOTED WITH RED AND BLUE, RESPECTIVELY.

Method	DukeMTMC-reID	
	r-1	mAP
BoW [9]	25.13	12.17
PAN [40]	71.59	51.51
FMN [41]	74.51	56.88
SVDNet [42]	76.7	56.8
HA-CNN [39]	80.5	63.8
Deep-person [43]	80.9	64.8
MLFN [38]	81.2	62.8
PCB [14]	83.3	69.2
Part-aligned[44]	84.4	69.3
US-GAN [11]	67.68	47.13
PN-GAN [26]	73.58	53.2
Liu et al. [24]	78.52	56.91
FD-GAN [27]	80.0	64.5
ResNet50-1 (Baseline-1)	74.48	58.95
ResNet50-2 (Baseline-2)	78.04	60.56
Ours	83.46	68.10
Ours (rerank)	84.01	69.94

TABLE V

COMPARISON OF RESULTS IN THE CUHK03 DATASET. THE BEST AND SECOND BEST RESULTS ARE DENOTED WITH RED AND BLUE, RESPECTIVELY.

Method	CUHK03			
	r-1	r-5	r-10	mAP
DeepReid [3]	19.89	50.0	64.0	-
LSTM Siamese [30]	57.3	80.1	88.3	-
PIE [5]	67.1	92.2	96.6	-
Gated Siamese [19]	68.1	88.1	94.6	-
PDC [34]	78.92	94.83	97.15	-
DLPAR [32]	81.6	97.3	98.56	-
SpindleNet [22]	88.5	97.8	98.6	-
SVDNet [42]	81.8	-	-	84.8
JLML [37]	83.2	98.0	99.4	-
US-GAN [11]	84.6	97.6	98.9	87.4
PN-GAN [26]	79.76	93.79	98.56	-
FD-GAN [27]	92.6	-	-	91.3
ResNet50-1 (Baseline-1)	83.30	88.54	92.13	76.51
ResNet50-2 (Baseline-2)	85.00	92.51	96.84	80.10
Ours	91.56	96.14	98.09	90.69
Ours (rerank)	92.24	97.35	98.86	91.14

4) *Failure Cases*: A few inconsistencies in the image generation part has been observed. Although our model generates qualitatively good images, the very fine details such as details of the face and the graphics of the t-shirt, are not properly reconstructed, which is primarily because of the low resolution of the input image.

VI. CONCLUSION

This paper proposes a novel person re-identification pipeline using GANs. The image generation capability of pose transformation GAN (pt-GAN) is used to model the subject under selected poses, which is defined by the Pose Clustering module,



Fig. 8. Comparison of retrieved images for 2 sample query image. The correct retrievals are denoted with green border. The Proposed model (ours) shows better accuracy compared to the two baseline models.

TABLE VI
COMPARISON OF RESULTS IN THE CUHK01 DATASET. THE BEST AND SECOND BEST RESULTS ARE DENOTED WITH RED AND BLUE, RESPECTIVELY.

Method	CUHK01		
	r-1	r-5	r-10
Ahmed et al. [2]	47.53	71.5	80.0
DeepRanking [45]	50.41	75.93	84.07
Ensembles [46]	53.4	76.3	84.4
ImpTrpLoss [18]	53.7	84.3	91.0
GOG [47]	57.8	79.1	86.2
Quadruplet [48]	62.55	83.44	89.71
NullReid [49]	69.09	86.87	91.77
PersonNet [50]	71.1	90.1	95.0
SpindleNet [22]	79.9	94.4	97.1
PN-GAN [26]	67.65	86.64	91.82
ResNet50-1 (Baseline-1)	64.89	83.76	89.84
ResNet50-2 (Baseline-2)	67.82	86.07	91.57
Ours	73.2	91.25	96.24
Ours (rerank)	73.95	93.04	96.97

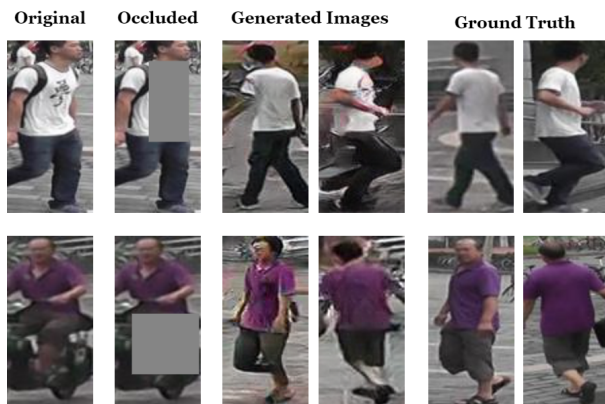


Fig. 9. Occlusion invariance in the proposed pt-GAN model. The occluded image is provided as the input to the model and the generated images are compared to the ground truth. The pt-GAN model is able to successfully reconstruct the occluded regions in various poses.

TABLE VII
COMPARISON OF RANK-1 ACCURACY WITH THE LATEST STATE-OF-THE-ART MODELS IN PERSON REID. THE IMPROVEMENT (%) OVER THE BASELINE HAVE BEEN TAKEN AS THE METRIC OF COMPARISON, WHERE THE PROPOSED MODEL EXCELS.

Dataset	Method	Rank-1		
		Baseline	Final	Improvement
Market-1501	Auto-ReID [51]	94.8	95.4	0.63
	pf-SAN [52]	93.8	94.7	0.96
	AdaptiveReID [53]	94.6	96	1.48
	VA-ReID [54]	94.7	96.79	2.21
	ABD-net [55]	91.5	95.6	4.48
	part-aligned [44]	88.8	93.4	5.18
	AlignedReid [56]	86.3	91.8	6.37
	DCDS [57]	87.5	94.1	7.54
	st-ReID [58]	91.2	98.1	7.57
	PCB [16]	86.7	93.8	8.19
	Ours	83.69	91.76	9.64
DukeMTMC-reID	AdaptiveReID [53]	88.00	92.20	4.77
	pf-SAN [52]	83.30	89.00	6.84
	VA-ReID [54]	87.39	93.85	7.39
	ABD-net [55]	82.80	89.00	7.49
	st-ReID [58]	83.80	94.40	12.65
	Ours	74.48	84.01	12.80
CUHK03	AlignedReid [56]	83.8	92.4	10.26
	DCDS [57]	87.7	95.8	9.24
	pf-SAN [52]	63.7	69.7	9.42
	Auto-ReID [51]	75	77.9	3.87
	Ours	83.3	92.24	10.73

and then features from all these representations are extracted and combined using the Feature Fusion Network. This end-to-end feature extraction model is used to extract a single viewpoint-invariant feature for each query, and the ranking is performed using feature similarity. Our model performs on par with the state-of-the-art models and outperforms the higher accuracy models in terms of improvement over baseline. We believe that our work will be useful to the community where unsupervised learning methods such as GANs can benefit person re-ID applications.

The possibility of further improving the overall re-ID performance using a different metric for matching i.e. retrieval is yet to be studied. The incorporation of attention maps or



Fig. 10. Comparative study of feature distances in the presence of rotation and illumination variation. The intra-class feature distances is significantly lower compared to that of the inter-class feature distance. The feature distances are scaled by a value of 10^3

part-based information in addition to the proposed method can also be a potential future work.

REFERENCES

- [1] A. Plantinga, "Things and persons," *The Review of Metaphysics*, pp. 493–519, 1961.
- [2] E. Ahmed, M. Jones, and T. K. Marks, "An improved deep learning architecture for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3908–3916.
- [3] W. Li, R. Zhao, T. Xiao, and X. Wang, "Deepreid: Deep filter pairing neural network for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 152–159.
- [4] J. Johnson, S. Yasugi, Y. Sugino, S. Pranata, and S. Shen, "Person re-identification with fusion of hand-crafted and deep pose-based body region features," *arXiv preprint arXiv:1803.10630*, 2018.
- [5] L. Zheng, Y. Huang, H. Lu, and Y. Yang, "Pose-invariant embedding for deep person re-identification," *IEEE Transactions on Image Processing*, vol. 28, no. 9, pp. 4500–4509, 2019.
- [6] M. Saquib Sarfraz, A. Schumann, A. Eberle, and R. Stiefelhagen, "A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 420–429.
- [7] S. Bak, F. Martins, and F. Bremond, "Person re-identification by pose priors," in *Image Processing: Algorithms and Systems XIII*, vol. 9399. International Society for Optics and Photonics, 2015, p. 93990H.
- [8] A. Karmakar and D. Mishra, "A robust pose transformational gan for pose guided person image synthesis," in *National Conference on Computer Vision, Pattern Recognition, Image Processing, and Graphics*. Springer, 2019, pp. 89–99.
- [9] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1116–1124.
- [10] N. Martinel and C. Micheloni, "Re-identify people in wide area camera network," in *2012 IEEE computer society conference on computer vision and pattern recognition workshops*. IEEE, 2012, pp. 31–36.
- [11] Z. Zheng, L. Zheng, and Y. Yang, "Unlabeled samples generated by gan improve the person re-identification baseline in vitro," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3754–3762.
- [12] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino, "Custom pictorial structures for re-identification," in *Bmvc*, vol. 1, no. 2, 2011, p. 6.
- [13] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani, "Person re-identification by symmetry-driven accumulation of local features," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 2010, pp. 2360–2367.
- [14] C. Weinrich, M. Volkhardt, and H.-M. Gross, "Appearance-based 3d upper-body pose estimation and person re-identification on mobile robots," in *2013 IEEE International Conference on Systems, Man, and Cybernetics*. IEEE, 2013, pp. 4384–4390.
- [15] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification: Past, present and future," *arXiv preprint arXiv:1610.02984*, 2016.
- [16] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 480–496.
- [17] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [18] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1335–1344.
- [19] R. R. Variator, M. Haloi, and G. Wang, "Gated siamese convolutional neural network architecture for human re-identification," in *European conference on computer vision*. Springer, 2016, pp. 791–808.
- [20] Y.-G. Lee, S.-C. Chen, J.-N. Hwang, and Y.-P. Hung, "An ensemble of invariant features for person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 3, pp. 470–483, 2016.
- [21] Y.-J. Cho and K.-J. Yoon, "Improving person re-identification via pose-aware multi-shot matching," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1354–1362.
- [22] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, and X. Tang, "Spindle net: Person re-identification with human body region guided feature decomposition and fusion," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1077–1085.
- [23] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [24] J. Liu, B. Ni, Y. Yan, P. Zhou, S. Cheng, and J. Hu, "Pose transferrable person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4099–4108.
- [25] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [26] X. Qian, Y. Fu, T. Xiang, W. Wang, J. Qiu, Y. Wu, Y.-G. Jiang, and X. Xue, "Pose-normalized image generation for person re-identification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 650–667.
- [27] 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 *Advances in Neural Information Processing Systems*, 2018, pp. 1222–1233.
- [28] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh, "Realtime multi-person 2d pose estimation using part affinity fields," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7291–7299.
- [29] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *arXiv preprint arXiv:1708.04896*, 2017.
- [30] R. R. Variator, B. Shuai, J. Lu, D. Xu, and G. Wang, "A siamese long short-term memory architecture for human re-identification," in *European conference on computer vision*. Springer, 2016, pp. 135–153.
- [31] X. Liu, H. Zhao, M. Tian, L. Sheng, J. Shao, S. Yi, J. Yan, and X. Wang, "Hydraplus-net: Attentive deep features for pedestrian analysis," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 350–359.
- [32] L. Zhao, X. Li, Y. Zhuang, and J. Wang, "Deeply-learned part-aligned representations for person re-identification," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 3219–3228.
- [33] S. Bai, X. Bai, and Q. Tian, "Scalable person re-identification on supervised smoothed manifold," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2530–2539.
- [34] C. Su, J. Li, S. Zhang, J. Xing, W. Gao, and Q. Tian, "Pose-driven deep convolutional model for person re-identification," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 3960–3969.

- [35] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [36] M. Geng, Y. Wang, T. Xiang, and Y. Tian, "Deep transfer learning for person re-identification," *arXiv preprint arXiv:1611.05244*, 2016.
- [37] W. Li, X. Zhu, and S. Gong, "Person re-identification by deep joint learning of multi-loss classification," *arXiv preprint arXiv:1705.04724*, 2017.
- [38] X. Chang, T. M. Hospedales, and T. Xiang, "Multi-level factorisation net for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2109–2118.
- [39] W. Li, X. Zhu, and S. Gong, "Harmonious attention network for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2285–2294.
- [40] Z. Zheng, L. Zheng, and Y. Yang, "Pedestrian alignment network for large-scale person re-identification," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 10, pp. 3037–3045, 2018.
- [41] G. Ding, S. Khan, Z. Tang, and F. Porikli, "Let features decide for themselves: Feature mask network for person re-identification," *arXiv preprint arXiv:1711.07155*, 2017.
- [42] Y. Sun, L. Zheng, W. Deng, and S. Wang, "Svdnet for pedestrian retrieval," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3800–3808.
- [43] X. Bai, M. Yang, T. Huang, Z. Dou, R. Yu, and Y. Xu, "Deep-person: Learning discriminative deep features for person re-identification," *Pattern Recognition*, vol. 98, p. 107036, 2020.
- [44] Y. Suh, J. Wang, S. Tang, T. Mei, and K. Mu Lee, "Part-aligned bilinear representations for person re-identification," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 402–419.
- [45] S.-Z. Chen, C.-C. Guo, and J.-H. Lai, "Deep ranking for person re-identification via joint representation learning," *IEEE Transactions on Image Processing*, vol. 25, no. 5, pp. 2353–2367, 2016.
- [46] S. Paisitkriangkrai, C. Shen, and A. Van Den Hengel, "Learning to rank in person re-identification with metric ensembles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1846–1855.
- [47] T. Matsukawa, T. Okabe, E. Suzuki, and Y. Sato, "Hierarchical gaussian descriptor for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1363–1372.
- [48] W. Chen, X. Chen, J. Zhang, and K. Huang, "Beyond triplet loss: a deep quadruplet network for person re-identification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 403–412.
- [49] L. Zhang, T. Xiang, and S. Gong, "Learning a discriminative null space for person re-identification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1239–1248.
- [50] L. Wu, C. Shen, and A. v. d. Hengel, "Personnet: Person re-identification with deep convolutional neural networks," *arXiv preprint arXiv:1601.07255*, 2016.
- [51] R. Quan, X. Dong, Y. Wu, L. Zhu, and Y. Yang, "Auto-reid: Searching for a part-aware convnet for person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 3750–3759.
- [52] H. Wang, Y. Fan, Z. Wang, L. Jiao, and B. Schiele, "Parameter-free spatial attention network for person re-identification," *arXiv preprint arXiv:1811.12150*, 2018.
- [53] X. Ni, L. Fang, and H. Huttunen, "Adaptivereid: Adaptive l2 regularization in person re-identification," *arXiv preprint arXiv:2007.07875*, 2020.
- [54] Z. Zhu, X. Jiang, F. Zheng, X. Guo, F. Huang, X. Sun, and W. Zheng, "Viewpoint aware loss with angular regularization for person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 13 114–13 121.
- [55] T. Chen, S. Ding, J. Xie, Y. Yuan, W. Chen, Y. Yang, Z. Ren, and Z. Wang, "Abd-net: Attentive but diverse person re-identification," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 8351–8361.
- [56] X. Zhang, H. Luo, X. Fan, W. Xiang, Y. Sun, Q. Xiao, W. Jiang, C. Zhang, and J. Sun, "Alignedreid: Surpassing human-level performance in person re-identification," *arXiv preprint arXiv:1711.08184*, 2017.
- [57] L. Tesfaye Alemu, M. Pelillo, and M. Shah, "Deep constrained dominant sets for person re-identification," *arXiv e-prints*, pp. arXiv–1904, 2019.
- [58] G. Wang, J. Lai, P. Huang, and X. Xie, "Spatial-temporal person re-identification," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 8933–8940.