

Person Re-identification: A Retrospective on Domain Specific Open Challenges and Future Trends

Asmat Zahra^a, Nazia Perwaiz^a, Muhammad Shahzad^b and Muhammad Moazam Fraz^{a,c,*}

^aNational University of Sciences and Technology (NUST), Islamabad, 44000, Pakistan

^bTechnical University of Munich, D-80333 Munich, Germany

^cThe Alan Turing Institute, British Library, 96 Euston Road, London NW1 2DB, United Kingdom

ARTICLE INFO

Keywords:

Visual Surveillance
Person Re-Identification
Literature Survey
Deep Learning
Open Challenges
Specific application-driven

ABSTRACT

Person re-identification (Re-ID) is one of the primary components of an automated visual surveillance system. It aims to automatically identify/search persons in a multi-camera network having non-overlapping field-of-views. Owing to its potential in various applications and research significance, a plethora of deep learning based re-ID approaches have been proposed in the recent years. However, there exist several vision related challenges, e.g., occlusion, pose scale & viewpoint variance, background clutter, person misalignment and cross-domain generalization across camera modalities, which makes the problem of re-ID still far from being solved. Majority of the proposed approaches directly or indirectly aim to solve one or multiple of these existing challenges. In this context, a comprehensive review of current re-ID approaches in solving these challenges is needed to analyze and focus on particular aspects for further advancements. At present, such a focused review does not exist and henceforth in this paper, we have presented a systematic challenge-specific literature survey of 230+ papers between the years of 2015-21. For the first time a survey of this type have been presented where the person re-ID approaches are reviewed in such solution-oriented perspective. Moreover, we have presented several diversified prominent developing trends in the respective research domain which will provide a visionary perspective regarding ongoing person re-ID research and eventually help to develop practical real world solutions.

1. Introduction


In recent years, person re-identification has received much interest owing to its widespread application prospects in numerous fields including intelligent video surveillance [1], robotics [2] and human-computer interaction [3] *etc.* Specifically, it is one of the fundamental components of an automated visual surveillance system where for public safety and security in a smart environment, an individual person may be automatically identified and tracked in videos (or images) acquired through multiple non-overlapping cameras installed on public places like airports, banks, cantonments, parks, streets, educational institutes *etc.* Since it is simply not feasible to rely on manual human intervention to identify a person of interest in huge amount of video data collected on daily basis, therefore a plethora of approaches have been proposed by vision researchers that aim to automate this highly challenging problem.

Methodologically, person re-id refers to identifying and tracking a person in multi-network non-overlapping cameras installed in indoor and outdoor environments. Given an image of a person captured from one camera, the task of person re-id is to identify this person from a pre-stored gallery set captured by other multiple cameras.

Despite growing trend in the number of publications appearing in top venues achieving increasingly higher accuracy on the existing benchmark datasets, the problem is still far from being solved to be translate into real world settings. This can be attributed to the number of challenges (e.g., occlusion, variations in person pose, viewpoint variations, misalignment, poor resolution *etc.*) that makes the problem extremely hard and needs to be resolved to bridge the performance gap between research (benchmark specific) and real-world environments.

1.1. Scope/Objective of the Review

In this paper, we have targeted the most popular challenges in person re-id to perform systematic challenge-wise review of the published approaches. In this context, we have collected papers from top conferences and journals for the years from 2015 to 2021. The progress in papers addressing each challenge and its influence on published results is

 moazam.fraz@seecs.edu.pk (M.M. Fraz)

ORCID(s): 0000-0003-0495-463X (M.M. Fraz)



Figure 1: Graphical view of enlisted challenges.(From left to right) (a) Occlusion, (b) Illumination variance, (c) Pose variance, (d) Background clutter, (e) Misalignment, (f) Scale difference, (g) Viewpoint variance and (h) Low Resolution

comprehensively reviewed. The specific challenges in re-id that are mainly considered in this review include: occlusion, pose variance, background clutter, misalignment, scale difference, illumination variance, viewpoint variance, low resolution and cross-domain or generalization. Particularly, the proposed review makes many-fold contributions. For instance, we provide in-depth analysis on impact of most popular re-id challenges by discussing the work on each challenge in top computer vision conferences and journals. This provides insights on complexities that arise due to each challenge in the whole re-id process. Moreover, based on reviewed progress, the best performing architectures achieving state-of-the-art (SOTA) results on each challenge (as shown in the Fig. 1) are highlighted and critically analyzed. Furthermore, we attempt to make future directions for researchers by comprehensively reviewing publications relevant to each challenge and discuss the limitations and benefits to lessen the gap between close-world and real-world implementations. Lastly, in addition to reporting trends and highlighting interesting approaches, we distil our analysis into few recommendations in the hope of fostering reproducible and efficient research in the field. For the readers from any of these scenario, this survey also present comprehensive information of how the challenges have been addressed in past and how various components of deep learning (DL) can be utilized to contribute in improving the person re-id considering the influence of each individual challenge.

1.2. Comparison with Previous Reviews

There already exist few review articles on person re-id [4, 5, 6, 7, 8]. Each one of them focus on different aspect of the re-id problem. For instance, in [4] both the hand-crafted and deep learning approaches based on image and video data have been reviewed. Similarly, advantages and disadvantages of the traditional and deep learning based approaches are critically analyzed in [5]. [8] focused on application-driven methods that are designed for specific applications and defined as generalized open-world re-id. In [6] six different learning methods including identification, verification, distance metric learning, part-based, video-based and data augmentation based deep models are comprehensively reviewed. Recently [7] person re-id is reviewed using open and closed world setting while keeping challenges of re-id perspectives in view. These different perspectives for close world setting includes feature representation, metric learning and ranking optimization. For open world setting, heterogeneity, end-to-end, semi or unsupervised learning, robust model learning with noisy data annotations and open-set person re-id are the considered perspectives.

All the aforementioned reviews have comprehensively provided the short-comings and benefits of considered methods and settings. Moreover, they have also provided the insightful future directions. However, none of them have systematically reviewed the influence of popular challenges on performance results and the role of datasets in resolving

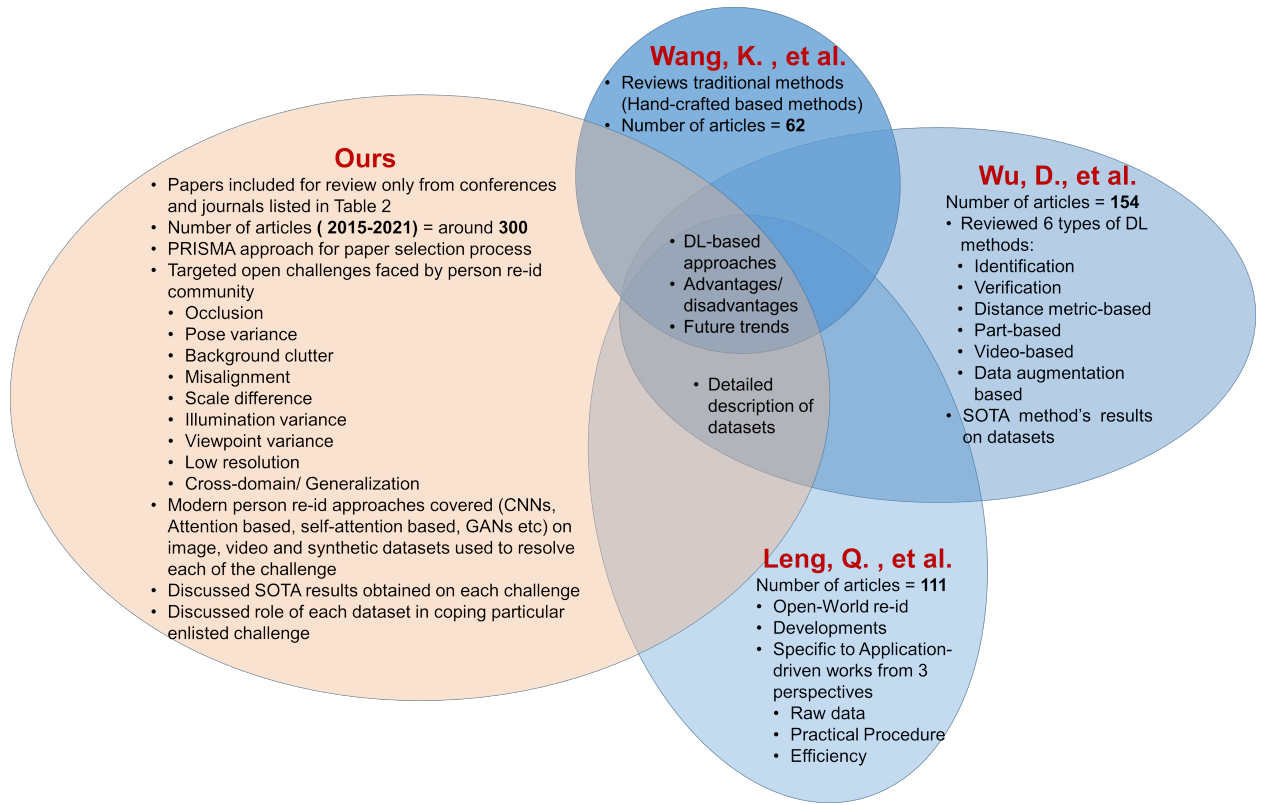


Figure 2: Comparison with recent published reviews on person re-id

these challenges. Specifically, no review exists that comprehensively addressed the challenges (mentioned in section 1) of the person re-id and their proposed DL based solutions. From 2015 to 2021 numerous articles have been published and each of these addressed the person re-id challenge in a specific way. Some of challenges are open to the re-id world and still not addressed properly. This motivated us to write this survey article which have comprehensively reviewed how all open challenges are addressed in past and how are results getting improved with respect to each challenge using DL based methods. The difference of our survey from the existing surveys can be visualized in Fig. 2.

2. Survey Methodology

This survey paper is organized in six sections, Section 1 introduces the domain, discusses the scope, objective and rationale of this paper, and highlights main contribution of this review by providing a detailed comparison with recent survey articles on person re-id. Section 2 illustrates the data collection methodology used for selection the articles include in this review. A comparative account on the publicly available datasets and performance metrics used to report result of person re-id methodologies is given in Section 3. Deep learning based approaches are used to solve person re-id challenges in the recent years. For the readers from inter disciplinary domain the deep learning approaches are briefly described in section 4. In Section 5, the open challenges in person re-id are discussed and the methodologies proposed to address these challenges are critically analyzed. Section 6 overviews the impact of challenges on results. Role of datasets in resolving the identified challenges along with limitation and benefits. And then it concludes the paper with possible future directions.

2.1. Study Selection

For study selection we used PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) [9]. Although PRISMA is used primarily in medical domain but because of its vast benefits we have used it for our review as well and obtained valuable results. Fig. 3 shows the summary of paper selection process. As per PRISMA guidelines,

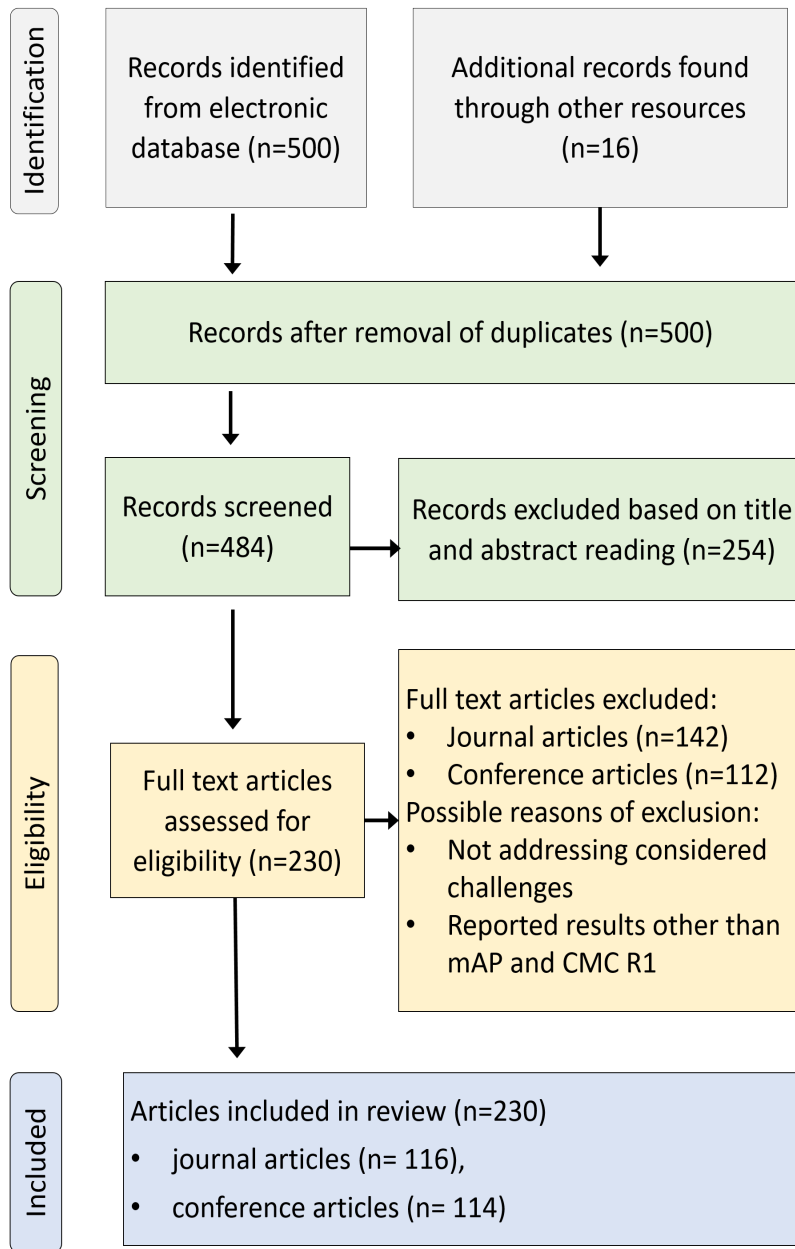


Figure 3: Summary of paper selection process (PRISMA).

our selection of articles for this review consists of two stages. In first stage all the articles that did not meet the eligibility criteria were excluded. While in the second stage we have studied full text reports to find the relevant articles. After shortlisting all the articles, we have excluded those paper which have reported results only through graphs and on such performance measure and datasets which were not so widely adopted by research community. In complex scenario where paper selection become difficult due to some ambiguity a discussion with senior member is organized to reach mutual final decision.

Table 1

Inclusion and exclusion criteria of selected articles

Inclusion criteria	Exclusion criteria
-Articles that address the identified challenges enlisted in table 2	-Articles in which qualitative evaluation of results are missing
-Provide detailed summary of proposed architecture including training parameters	-Papers that does not address the identified challenges enlisted in table 2
-Articles that are based on deep learning techniques	-Survey papers
-Articles that are published in journals and conferences enlisted in table 3 between January 2015 to October 2021	-Reported results on metrics other than Rank-1 and CMC
	-Papers that have used datasets that were not so widely adopted by research community

2.2. Data Extraction Methods

We have collected an initial list of articles from Springer, Google Scholar, IEEE Xplore and Elsevier. We have considered top seven journals and three conferences held from 2015 to 2021 for review. We have used the following terms (or matching to these) to search relevant articles:

- (a) Person Re-Identification
- (b) Deep learning
- (c) Supervised person re-id
- (d) Semi-supervised person re-id
- (e) Unsupervised person re-id
- (f) Pose variations
- (g) Body misalignment
- (h) Attention based approach
- (i) Camera View(Viewpoint)
- (j) End-to-end learning

Search results are further enhanced by combining the mentioned terms using logical operators in a way like: ('a') AND ('b' OR 'c' OR 'd' OR 'e' OR 'f' OR 'g' OR 'h' OR 'i' OR 'j'). The papers were then excluded or included based on the criteria listed in table 1. We first read the title of article for final selection. In case title does not clearly fall in our inclusion and exclusion criteria then we also read the abstract and conclusion section of the article as well. After that we start full reading for data collection. And articles that does not match with our criteria were not include in this review.

2.3. Data Synthesis

To make our review more useful in a sense that other researcher can contribute into it in future for the purpose for extension in review in multiple perspective, a data extraction sheet was developed that describe the multiple rel event data items to be extracted from the articles. Around 30 data items were used for metadata extraction from each article. These data items were classified in seven categories: Origin of the article, The challenges addressed, Details on the proposed methodology, Implementation detail, Reproducibility and code availability, The performance metrics and the Datasets used. Table 2 shows the category wise distribution and description of each of the data item. The results will

Table 2

Summary of data-items extracted from each article.

Category	Data Item & Description
Origin of article	Conference/ Journal Name Publication year & venue Article Title Image/Video dataset or both CNN, ViT, GCN, GAN based approaches used in the paper
Challenges	Occlusion : Blockage or hiding of target person in image Pose Variation : Particular person appeared in different positions Background clutter : A pattern present in background resembles with pattern of person's wearing in image Body Misalignment : Person not aligned according to viewing angle in image Scale : Person appearing in different sizes in an image Illumination : Light variations in an image Viewpoint variation : Change in position of capturing camera Resolution : Clarity/detail in input image Cross-domain /Generalization : Images belong to multiple domains
Proposed Methodology	Description of proposed method in article
Implementation Details	Implementation Framework & Platform used Base Model : DL baseline model used as backbone Training approach of model Take single/multiple images as query Data-set used in pre-trained model Batch Normalization Scheme used Batch size considered for training Type of pooling used in the article Learning rate-decay Data Augmentation technique used
Reproducible	Code Availability for public use or not
Results	Performance Metric used for reporting the results
Data-sets	Name(s) of data-sets used

be stored in a spreadsheet which will be made public for interested researchers intended to extract more information or analyse a different perspective.

The first category extracts the details about the origin of the article whether the article is published in a conference or a journal, the year of publishing, the title of article, which problem the particular article is addressing, image based or video based and finally which methodology is used i.e. CNN, Vision Transformer (ViT), GANs or Graph CNNs etc. In the second category, the details of enlisted person re-id challenge (Poor resolution, background *etc.*) addressed by the article is explored. In the third category, a brief extract of the proposed methodology in the article is analysed. The fourth category comprehensively extracts the implementation detail including implementation framework, deep learning based methodology, base model, training approach, batch normalization etc. The fifth category tells about public availability of implementation. The sixth category gathers the information about the quantitative performance measures of proposed work. Finally, the details about the datasets used for evaluation is reported in the last category. The Taxonomy of our review is shown in Fig. 4.

2.4. Results

There were 516 articles collected using the selection criteria. After removing irrelevant articles database contain 500 articles. According to our inclusion and exclusion criteria more 254 articles were excluded and therefore we

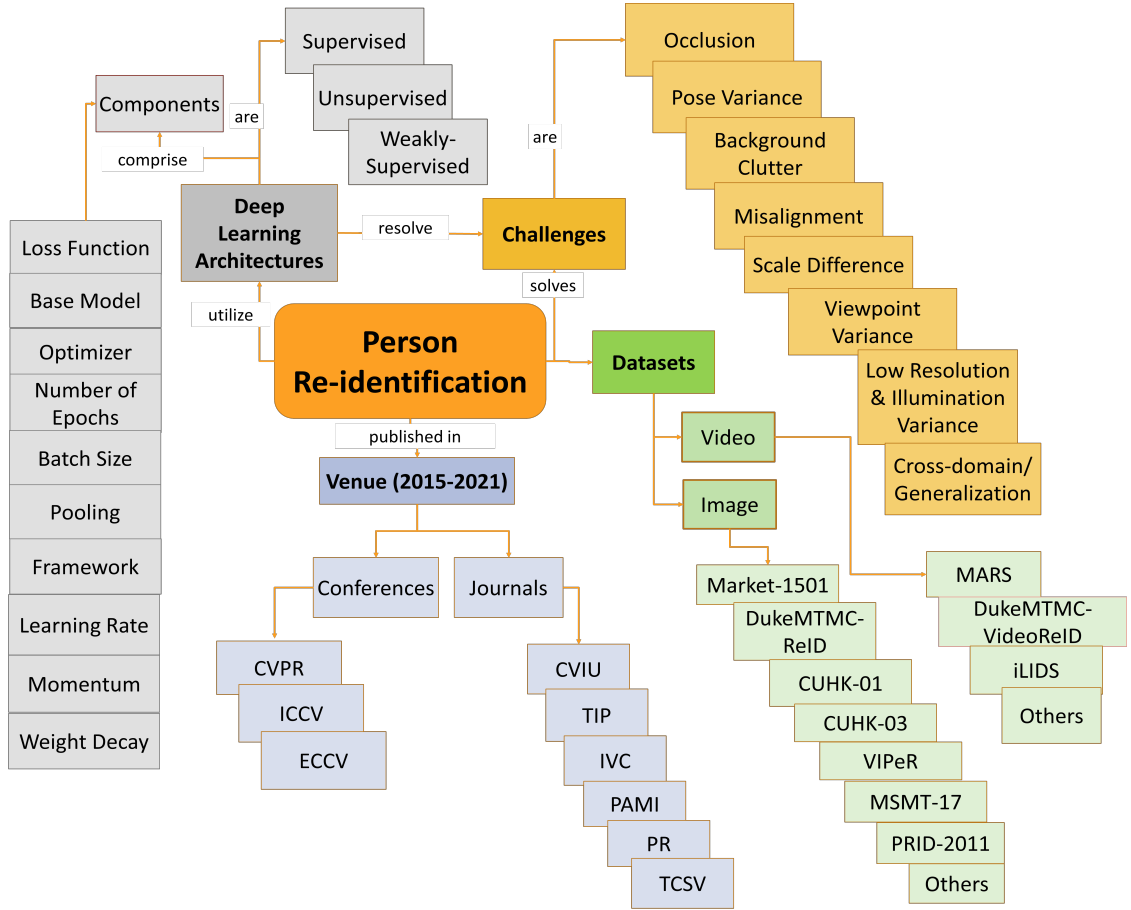


Figure 4: Taxonomy of the Review

selected 230 articles (116 Journal & 114 Conference articles) for review. The table 3 depicts the publication venue wise summary of the articles selected for this review, and Fig. 5 illustrates the growing trend of person re-id publications each year in reputed venues.

3. Datasets and performance evaluation

Several image based and video based datasets for person re-identification models has been released. In past few years several review papers review the available datasets. But they have not reported the progress on each challenge against each dataset. We review several image and video based datasets by reporting the progress of each challenge on each dataset and how dataset support the challenge to be addressed effectively. The attributes of each of the public datasets is summarized in the table 4.

3.1. Performance Evaluation Metrics

For evaluation purpose, person re-id models use two widely known measurements named Cumulative Matching Characteristics (CMC) also known as rank-k accuracy and mean Average Precision (mAP). CMC represents the probability that accurate match of image appears in the top-k ranked retrieved results. CMC will be considered accurate when only one ground-truth exists for each query, since it only considers the first match in evaluation process. However, the gallery set usually contains multiple ground truths in a large camera network, and CMC cannot completely reflect the discriminability of a model across multiple cameras. The other widely used metric, mAP measures the average retrieval performance with multiple ground- truths. It is originally widely used in image retrieval process.

Table 3
Summary of articles collected

Sr.No	Venues (2015-2021)	Found	Filtered	Reviewed
1	Computer Vision and Pattern Recognition (CVPR)	131	72	59
2	IEEE International Conference on Computer Vision (ICCV)	50	18	32
3	European Conference on Computer Vision (ECCV)	45	22	23
4	Elsevier Computer Vision and Image Understanding (CVIU)	6	0	6
5	IEEE Transactions on Image Processing (TIP)	91	46	45
6	Image and Vision Computing (IVC)	8	3	5
7	IEEE transactions on pattern analysis and machine intelligence (PAMI)	42	27	15
8	International Journal of Computer Vision (IJCV)	8	4	4
9	Pattern Recognition (PR)	42	27	15
10	IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)	51	31	20
	Total	484	254	230

3.2. Image based datasets

Several image based person re-id datasets are available. We have considered top 6 image based person re-id datasets *i.e.* Market-1501, DukeMTMC-ReID, CUHK03, CUHK01, VIPeR and MSMT-17 for our review on which person re-id models achieved SOTA performance. Major person re-id challenges are addressed by using these mentioned datasets. Summary of number of papers on major challenges is shown in the Fig. 6.

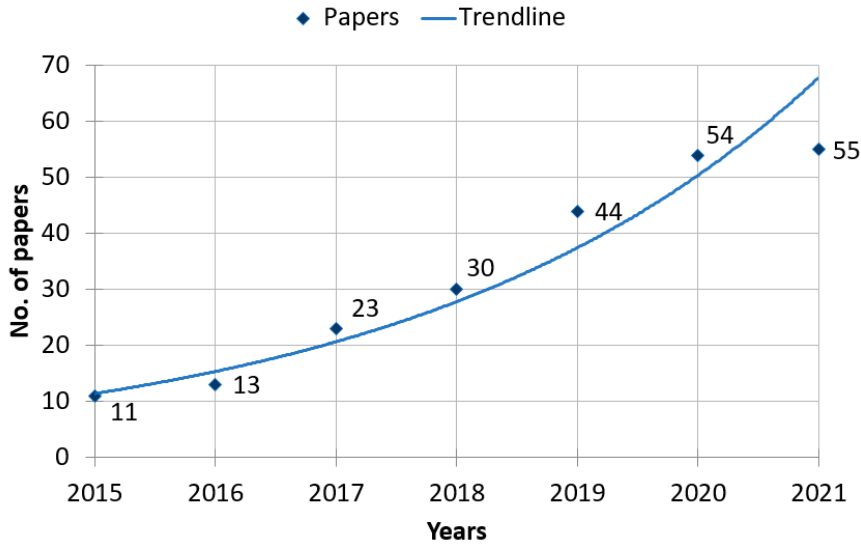


Figure 5: The yearly frequency of selected re-id papers for this review

3.2.1. Market-1501

The Market-1501 [10] dataset is specifically for person re-id and proposed in 2015. It was collected in front of a supermarket in Tsinghua University. A total of six cameras were used, including five high-resolution cameras, and one low resolution camera. Overlap exists among different cameras. Overall, this dataset contains 32,668 annotated bounding boxes of 1,501 identities. Each annotated identity is present in at least two cameras, so that cross-camera search can be performed. Market-1501 is one of the most trending and widely used dataset by researchers to address multiple person re-id challenges like illumination, viewpoint, pose variation and body misalignment. In 2019, this dataset helped to boost the performance of person re-id models. Market-1501 dataset provides the viewpoint angle and this property help the deep learning models in learning the deep features to resolve the viewpoint and pose variation challenge. In last five years many person re-id algorithms achieved the SOTA results by using Market-1501 dataset.

3.2.2. DukeMTMC-ReID

DukeMTMC-ReID [11] has significant potential, it provides access to details like frame level, ground truth, full frames and calibration information *etc.*). It correspond to images of 1,852 people existing across all the 8 cameras. It covers 1,413 unique identities with 22,515 bounding boxes that appear in more than one cameras. It also consists of 439 distractor identities with 2,195 bounding boxes that appear in only one camera. The size of the bounding box varies from 72*34 pixels to 415*188 pixels.

3.2.3. CUHK01

The image quality of CUHK01 [12] datasets is relatively good and this benefits the person re-id models to achieve good results and perform well in real world scenarios. This dataset was published in 2012 and consists of 3884 images of 971 people. Two disjoint cameras were used to capture different views. Each camera captures two images for each person, total of four images of a person.

3.2.4. CUHK03

CUHK03 [13] dataset is one of the largest dataset in 2014 and proved good for person re-id and deep learning models to report effective results. Dataset comprised of 13,164 images of 1,360 people. Images are captured by using six cameras. Each identity appears in two disjoint camera views (*i.e.* in each view there are 4.8 images on average). In CUHK03 bounding boxes are manually labeled and detected from Deformable Part Models (DPM).

3.2.5. VIPeR

VIPeR [14] is one the most challenging dataset, several researchers tested it and reported very widely and interesting results specially by addressing viewpoint variation challenge. As this dataset is one of the oldest dataset and still in

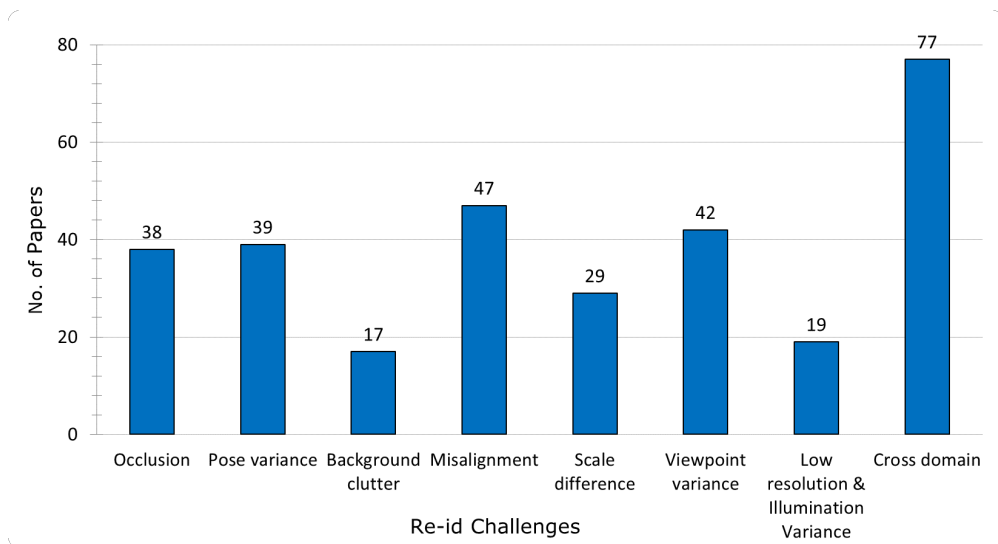


Figure 6: Count of papers on each enlisted challenge

Table 4

Properties of each dataset.

Sr.No	Dataset	Year	Environment	Identities	Cameras	Resolution	Label	BBoxes	Challenging Attributes
1	VIPeR [14]	2007	Campus	632	2	48×128	Hand	1,264	VV, IV
2	PRID-2011 [15]	2011	Outdoor	200	2	64×128	DPM/ GMMCP/ Hand	40,000	IV, VV,,BC
3	CUHK01 [12]	2012	Campus	971	2	60×160	Hand	3,884	VV,OCC
4	CUHK03 [13]	2014	Campus	1,360	10	Vary	DPM/Hand	13,164	VV,OCC
5	iLIDS-Vid [16]	2014	Airport	300	2	Vary	Hand	42,495	VV,IV,BC,OCC
6	Market-1501 [10]	2015	Campus	1,501	6	64×128	DPM/Hand	32,688	VV, PV, RES
7	MARS [17]	2016	Campus	1,261	6	256×128	DPM/ GMMCP	1,067,516	PV,IV,RES
8	DukeMTMC-RelD [11]	2017	Campus	1,404	8	Vary	Doppia/ Hand	36,411	VV,IV,BC,OCC
9	DukeMTMC-Video RelD	2017	Campus	1,404	8	Vary	Hand	36,411	VV,IV,BC,OCC
10	MSMT-17 [18]	2018	Campus	4,101	15	Vary	Faster RCNN	12,6441	VV,IV

¹ Viewpoint Variation² Illumination Variation ³ Background Clutter ⁴ Occlusion ⁵ Pose Variation ⁶ Resolution ⁷ Deformable Part Model [19] (A Pedestrian detector) ⁸ GMMCP- Generalized Maximum Multi Clique problem [20] (A tracker)

trending for person re-identification models. It contains around 632 identities and images are captured by two cameras one image per person. VIPeR also facilities with the viewpoint angle of each image. In our review we have observed that the in last five years most journal paper on person re-identification test their model on VIPeR dataset to address pose variation and viewpoint challenge.

3.2.6. MSMT-17

MSMT-17 [18] is one of the largest image based dataset containing 126441 images and 4101 identities. Images are captured in morning, noon and afternoon in campus. In our review we have observed that MSMT-17 is widely used dataset, although it contains similar viewpoint with Market-1501 dataset but this data commonly known for capturing the most complicated scenarios that's why several and almost all recent person re-id models test their models on this dataset and reported SOTA results. It mostly adopted to address pose variation, viewpoint and body misalignment challenges.

3.3. Video based datasets

Several video based person re-id datasets are available. We have considered top 3 video based person re-id dataset MARS, DukeMTMC-Video re-id and iLIDS for our review on which person re-id models achieved best performance. Major person re-id challenges are also addressed by using these video based datasets.

3.3.1. MARS (Motion Analysis and Re-identification Set)

MARS [17] dataset is an extended version of Market-1501 dataset. It is largest video based dataset having 1191003 images and with maximum crop size 256*128 among all other video based datasets. In our review we have found that MARS is specifically used to deal with pose variation, viewpoint and similarity measure challenges. Paper published in almost all top journals based on video based person re-id very few of them used this dataset to test their models

but these effective results explain that the MARS is significant dataset because all the tracklets and bounding box are generated automatically. This automatic generation makes learning faster.

3.3.2. iLIDS-VID

iLIDS [16] contain 42495 images with 300 identities. iLIDS dataset contain heavy occlusion in the captured images that's why it mostly used in person re-id models which address occlusion challenge. In our survey we have found that iLIDS majorly help to achieve good results in pose variation, viewpoint and similarity measure challenges. Therefore, in last five years very few journal papers that address these challenges used iLIDS dataset.

4. Deep Learning Approaches

In recent years, neural network-based deep learning algorithms become a popular branch of machine learning. Deep learning algorithms attempt to model high-level abstractions in data by using multiple processing layers with complex structures and are able to outperform state-of-the-art methods in many tasks in the fields of computer vision, natural language processing, robotics, etc. For the interdisciplinary readers, in this section we will present an overview of deep learning architectures used to address various challenges of person re-identification

4.1. Convolutional Neural Network

For the computer vision task, the input data are images of sizes usually ranging from several hundreds to several tens of thousands of pixels. If a neural network processed this input matrix with only fully connected neurons, the number of parameters to train would be very large, leading to a high risk of overfitting. The convolutional layer was invented to deal with this problem and became the first successfully trained deep neural network. The convolutional layer uses two basic ideas: local receptive fields and shared weights to reducing the complexity of the neural network. A CNN receives a matrix as the input, but connects a hidden node to only a small region of nodes in the input layer, since the spatial correlation is local. This region is called the local receptive field for the hidden node. Moreover all the mappings between a local receptive field and a hidden node share the same weights, since the features are not specific to some regions in an image. These two properties reduce dramatically the numbers of parameters of the network. By applying these two principals, a fully-connected neuron layer is transformed into a convolutional layer as follows:

$$y = \alpha(W * X + b) \quad (1)$$

where X is the input image W is the weight matrix, also called a 'filter' or 'kernel'. b is the bias term. The function α is an activation function and the operator $*$ represents the discrete convolution operation. For a two-dimensional image X as input, the convolution operator is defined as:

$$(W * X)(i, j) = \sum_m \sum_n X(m, n)W(i - m, j - n) \quad (2)$$

Intuitively, the output of the convolutional layer is formed by sliding the weight matrix over the image and computing the dot product. The resulting matrix is called "activation map" or "feature map". In image processing, convolution operations can be employed for edge detection, image sharpening and blurring just by using different the numeric values of the filter matrix. This means that different filters can detect different features from an image and capture the local dependencies in the original image. In convolutional layers, the convolutional kernel or filter, i.e. the coefficients of the weight matrix W , are learnt automatically by the backprop algorithm, and one layer usually contains several such convolution kernels and resulting feature maps.

In a CNN architecture, the image is passed through a series of convolutional, nonlinear, pooling layers and fully connected layers, and then generates the output. The Convolution layer is always the first layer where the image is entered. The reading of the input matrix begins at the top left of image with the help of a smaller matrix, which is called a filter. The filter produces convolution, i.e. moves along the input image. The filter multiplies its values by the original pixel values. All these multiplications are summed up. One number is obtained in the end. Since the filter has read the image only in the upper left corner, it moves further and further right by 1 unit performing a similar operation. After passing the filter across all positions, a matrix is obtained, but smaller then a input matrix. This operation, from a

human perspective, is analogous to identifying boundaries and simple colours on the image. But in order to recognize the properties of a higher level such as the trunk or large ears the whole network is needed.

The network will consist of several convolutional networks mixed with nonlinear and pooling layers. When the image passes through one convolution layer, the output of the first layer becomes the input for the second layer. And this happens with every further convolutional layer. The nonlinear layer is added after each convolution operation. It has an activation function, which brings nonlinear property. Without this property a network would not be sufficiently intense and will not be able to model the response variable (as a class label).

The pooling layer follows the nonlinear layer. It works with width and height of the image and performs a downsampling operation on them. As a result the image volume is reduced. This means that if some features (as for example boundaries) have already been identified in the previous convolution operation, than a detailed image is no longer needed for further processing, and it is compressed to less detailed pictures.

After completion of series of convolutional, nonlinear and pooling layers, it is necessary to attach a fully connected layer. This layer takes the output information from convolutional networks. Attaching a fully connected layer to the end of the network results in an N dimensional vector, where N is the number of classes from which the model selects the desired class.

Initially, the deep learning based re-id research utilized the ImageNet pre-trained weights to fine tune the most popular deep architecture ResNet on the person re-id or pedestrian datasets [21], [22].

4.1.1. Global Features Learning

In the past few years, with the rise of deep learning for image classification tasks, the development of various deep architectures and the availability of their pre-trained weight, various classification loss functions are proposed for deep learning base person representations [23] [24]. Furthermore self-learning of person features is endorsed very first time by using pairwise Siamese loss [13] and triplet loss [25] where features selection brings similar identities closer and different identities apart. Person attributes and the salient features in a person image are also worked out to focus on salient regions of a person images [26] [27] [28] [29] [30]. A technique for hard-identity mining is proposed to learn global features in an efficient way by Ristani et al in [31], where a weighted triplet loss function is introduced and a robust two streams metric learning model is proposed for person tracking and person re-id in parallel. In [32] a restraint and relaxation iteration (RRI) training scheme is used to propose SVDNet for person re-id where an Eigenlayer is introduced to find the correlation among the global features before fully connected layer of the deep model. However the spatial information loss of local cues need to be handled effectively for improvement of re-id performance.

4.1.2. Local Features Learning

In addition to learning global visual descriptors for person representation, the re-id research deals with the local vli2014deepreidusal cues as well. A famous strategy for working on the local/ part based person features is to divide the image in multiple local parts or strips. In the re-id approaches [13][33][34][35][36] the person cropped image is divided into various horizontal strips to emphasize on the local discriminative regions in an image. The methodology proposed in [13] divides the person image into three fixed horizontal parts and assumes that the top part of the image contains head, middle one contains torso and the bottom one contains legs. [33], [34], [37], [38] use another fixed number of local parts for local features whereas in [35] a divide and fuse strategy is adopted to focus on local parts of image, however the association among local parts of images is not catered for in any of these works.

One of the major limitation of all these convolution based attention networks is that these networks learn the dependency among immediate neighborhoods both at the initial layers of the network and at the deeper layers as well [39]. The structure and working mechanism of CNNs do not exhibit the learning of attention dependency at distant positions of an image or feature maps at a particular level. This arises the need to learn self-attention or intra-attention at each learning level (layer). The self-attention learns the associations among different parts of images and embed this information in the global representations.

4.2. Attention Based Deep Network

In the convolution based deep architectures, the attention regions of the images are captured at small receptive fields of the input image and are propagated towards deeper layers of the deep networks to have their aggregated perception [40] [41] [42] [43] [44]. This cumulative attention information provides good intuition about the regions of the images which can play an important role for the task of person re-id. But the deeper layers of CNNs are bound to view whatever the initial layers of a deep architecture have forwarded to them hence losing the global context to greater extent.

Over the last couple of years, with the advancement of hardware architectures and the availability of large-scale datasets, the attention based deep architectures are being explored extensively for the person re-identification ongoing research [45][46]. Attention mechanisms get intuition from the human vision system and focus on the most attentive regions of the image for decision making and suppress the less attentive regions.

The attention operations map the queries (Q) and a set of key-values (K, V) to the output vector. The queries, keys and values are the matrices and make the input for an image for which attention is to be computed. Two common kinds of attention functions are the multiplicative attention function i.e. scaled dot product attention, and the additive attention function. The most common method is the multiplicative function to compute the computationally efficient attention matrices instead of the additive attention method. The additive attention function is computationally expensive and uses a feed-forward network with a single hidden layer i.e 3 layer network (input + hidden + output), which involves one multiplication of the input vector by a matrix, then by another matrix, and then the computation of resultant vector. In contrast, the smart implementation of the scaled dot product attention computation does not break out the whole matrix multiplication algorithm and basically is a tight, easily parallelized operation.

By using variants of CNNs as a backbone, several re-id networks are designed that explicitly embed higher level attention in addition to the local attentions learnt by CNNs base architectures. Most of these attention learning networks opt multi-stream structures to learn the regional/ spatial attention along with deep convolutional representations [45][46]. In addition of learning the spatial attentions, the attentive but diverse re-id model [47] and critical attention learning mechanism [48] also focus to learn channel wise attentions in order to extract the most significant channels only as all the channels do not contain significant information. In addition to attentive channel information, [47] learns correlation among the attentive channels as well.

Another multi-task and attention aware learning network [49] proposed by Chen et al. demonstrated an explicit holistic attention branch to learn global attention in addition to a partial attention learning branch to learn local attention, however the network has additional requirements of the key-points for its local attention branch and learning self-attention within sparse parts of the image.

The above mentioned methods learn hard-level attention to perform person re-id, but the computation of only hard attention makes these networks less generalized. Harmonious attention networks [50] handle multi-level attention i.e. both of the hard attention (regional salience) and the soft attention (pixel level salience). It retains the generalization of CNNs as well but it still fails to capture the relationships among distant attentive regions within an image. All these multi-stream attention learning mechanisms need huge computational resources. Moreover, these networks do not learn the associations among far-distant attention regions of an image which are apparently dispersed all over the image but can provide better intuition to re-identify a person.

4.3. Self-Attention Based Deep Network

One of the major limitation of all these convolution based attention networks is that these networks learn the dependency among immediate neighborhoods both at the initial layers of the network and at the deeper layers as well [39]. The structure and working mechanism of CNNs do not exhibit the learning of attention dependency at distant positions of an image or feature maps at a particular level. This arises the need to learn self-attention or intra-attention at each learning level (layer).

The self-attention learns the associations among different parts of images and embed this information in the global representations. Transformer is a deep learning model for learning self-attentions among a given sequence of inputs and is currently state-of-the-art for solving the language problems [39]. However due to the large number of the pixels in images, it is not practical to use a standard transformer in its default mode for the images based classification tasks. A standard transformer encoder intakes a sequence of one dimensional data, learns self-attentions among the given sequences and embeds this self-attention information into the final representations of the data. The self attention computation (SA) is given in the equation 3.

$$SA(Q, K, V) = \left(\frac{\exp(QK^T)}{\sum_j \exp(QK^T)_j} * \frac{1}{\sqrt{d_k}} \right) V \quad (3)$$

In these methods, multiple parallel self-attention heads are used to run multiple parallel attention functions. The main objective is to learn attention from different positions and representation spaces and jointly update learnable parameters of each attention head as shown in equation 5. This mechanism efficiently learns the self-attention and associations

among local parts at distant positions in an image. The weights of the attention are set on the basis of pairwise similarity among the patches' sequences. These learnt self-attentions are further refined through the depth of transformer layers.

$$h_i = SA_{mul}(QW_i^Q, KW_i^K, VW_i^V) \quad (4)$$

where, h_i is the attention computed by i th head with resultant trainable parameter matrices i.e. W_i^Q, W_i^K, W_i^V , as shown in equation 4

$$MSA = F_x \left(\sum_i h_i \right) W^o \quad (5)$$

Finally, the multi-head self-attention module integrates the attention computed by each head using the function (F_x) of concatenation. It limits the information loss through a residual connection i.e. W^o , which contains the normalized output vector of the preceding layer. The W^o is integrated with a layer's output before submitting it to the next modules. The attention learning mechanism aims to learn the attention and its interdependence in a global manner.

For vision tasks, one potential way is to take the pixels sequences as one dimensional input data. However, typically the large number of pixels in images makes the self-attention computation across the whole set of pixels computationally very costly and thus consequently limits the use of transformers particularly for vision applications. Recently, [51] proposes an effectual pre-processing pipeline to handle the massive pixels of images data, which enables the use of transformers for vision based problems.

Luo et al. first time use the transformer based network for person re-id. [52] integrates the deep convolutions based re-id module with the pair-wise spatial transformer networks (STN) module to perform the partial person re-id. The spatial transformer networks module samples an affined image (a semantically corresponding patch) from the holistic image to match the partial image. The re-id module learns the embedding of holistic, partial and affined images, the STN module performance is influenced by the re-id module.

Taking inspiration from the state-of-the-art self-attention methods for natural language problems solution, the researchers explored the self-attention impact for person re-id task. In [53], the multi-scale convolutions are applied to the entire image and to the pre-defined three local parts of the image, i.e. upper, middle and bottom. The latent parts localization is performed by using spatial transformer networks to learn the self-attention. However, this work needs the positions of local parts and the value range constraint on the scale parameter as prerequisites.

5. Person Re-id Challenges and State-of-the-Art Results

In recent years, person re-id has gained impact full attention in the community of intelligent system and computer vision for various decisive applications. Although person re-id is comprehensively studied by researchers globally but still it is a challenging issue. When images were captured by non-over-lapping cameras under dynamic-environment are of low quality and in some images face or some other important features are not covered comprehensively. Conventional methods based on hand crafted algorithms and small-scale evaluation are not feasible because of their limited applications. In past, results based on subtracting background from frame to frame in multi-camera tracking are not enough due to variations in viewpoint, domain and illumination *etc.* However, the reliable re-id mainly involves the accurate response that is near to real time. This requires the availability or selection of good quality images to cope the challenges. Therefore, bunch of challenges are still not fully addressed *i.e.* occlusion, pose variations, background clutter, misalignment, scale, illumination, viewpoint changes, poor resolution and cross-domain or generalization.

In our review we have provided detailed progress on each of the challenges in last six years as shown in Fig. 6. How the results are getting better yearly and which datasets are mostly used to solve these challenges. How deep learning techniques have improved results on each challenge.

We have categorized the papers of each challenge on the basis of adopted algorithmic techniques *i.e.* CNN, attention and self-attention based approaches.

- **CNN based Re-id Solutions:** CNNs are widely used because of their in-depth learning. They consists of number of convolutional layers. Mainly used for image processing, segmenation, classification and other correlated data. A sliding filter is used to convolve over the entire image to gradually learn the portions of the image and its surroundings.

- **Attention based Re-id solutions:** Attention-based approaches are also in wide use due to their focusing behaviour of learning specific attributes. Their use enhances the focus on important part of data while ignoring the unused background information.
- **Self-Attention based Re-id Solutions:** While transformers are newly introduced in vision processing tasks to boost their performance at next level. They use multi-head self-attention mechanism to learn image embedding. Transformers can also be used in conjunction with CNNs however pure transformers applied directly to image patches can also produce best results.

5.1. Occlusion

Occlusion is caused by any overlapping object that may lead to wrong results. Persons are occluded by various environmental objects (traffic sign board, trees in parks, vehicles in parking) or by other pedestrians in person retrieval scenarios and make it difficult to track the movement of people. Visual illustration of occlusion is illustrated in the Fig. 1.

When a person is occluded, the features extracted from the whole image may contain the distracted information and leads to wrong results if the model is not able to differentiate between the occluded region and person region. Recently [25, 54, 55] few of them solved the occlusion challenge but did not consider the situation where the person is occluded due to different obstacles like cars, walls, shelves, poles and other people. In previous works to reduce the effect of obstructions, the occluded target person in the probe images is cropped manually and then use the non-occluded part as the new probe image. There are limitations in such approaches because it require more processing time due to manual cropping.

Some of the works on occlusion [54, 56] have achieved better performance by using part based models via part to part matching. It's not ideal to discard the occluded part from the target image, recently attention mechanism has been introduced in person re-id models which pay more attention to the non-occluded part during feature construction and preserve the most discriminating and informative appearance information and leads to better person retrieval results. In the literature there were many approaches reported to handle the occlusion. But we have divided the collected papers in three sub-categories: CNN-based, Attention-based and Transformer-based.

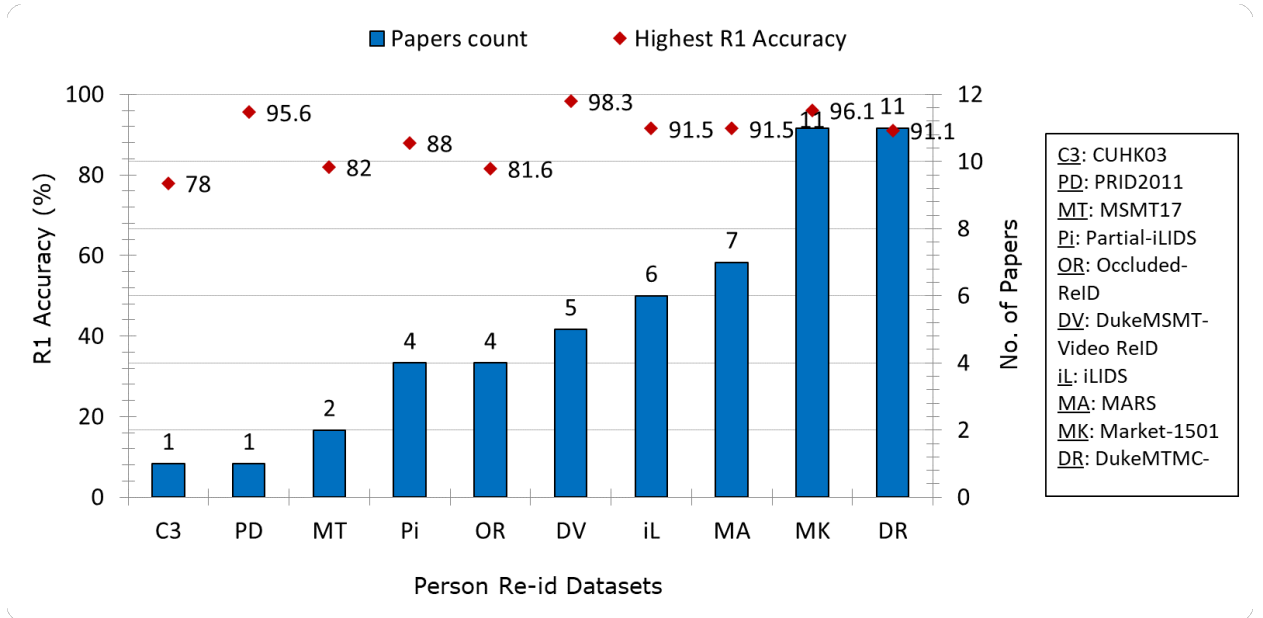


Figure 7: Progress on the challenge of occlusion for person re-id benchmarks

5.1.1. CNN-based approaches:

In [57] crowded scenes are challenged. For handling partial occlusion a simple occlusion-aware approach is presented. Fully convolutional network was used to obtain spatial features. In order to make the features more

discriminative a salience heat-map is created using combination of mask-guided and pose-guided layer. Saliency heat-maps are then used to guide the adaptive spatial matching. Adaptive matching is for assigning the larger weights to foreground human parts, so to obtain effective results as compared to existing state of the art approaches. Lingxiao He *et al.* [56] designed an occlusion and alignment free framework to obtain spatial pyramid features at multiple levels/scales. Fully convolution network was used to generate discriminative feature maps. No prior information about alignment was used. A CNN based estimation model was used to extract the semantic information in [58]. An adaptive graph based direction providing layer is introduced with a purpose to pass the information of relation to the nodes. This layer also stops the passing of irrelevant information by analysing how strong the link of nodes is, which is determined by amount of shared information. Moreover Cross graph layer was added to embed and learn the topology and alignment information of group of local features using graph matching technique.

A new large scale occluded person re-id dataset was introduced in [59]. Proposed framework was tested not only on existing benchmark datasets but also on the newly proposed occluded dataset. In the presented framework single-scale discriminative features are learnt without using any auxiliary module. A bounded distance loss makes the approach unique as this enables the model to learn discriminative features explicitly from occlusion-based augmented data. Jnrui Yang *et al.* in [60] a simple and effective mechanism was presented to elevate the occlusion challenge. Pose information was discretized to the level of visibility. In this manner the impact of occluded body regions is suppressed that has eventually helped in learning the robust and effective pose information on occluded person re-id.

Cairong Zhao *et al.* in [61] presented an Incremental Generative Occlusion Adversarial Suppression (IGOAS) network. The proposed model consists of two blocks *i.e.* a generative block and a global block. Former block generates the easy-to-hard occluded data samples to make the model robust to occlusion challenge while the later block usually extract the global features, suppresses the impact of occluded part and strengthens the focus on non-occluded part of body regions. Results are then concatenated to obtain the discriminative feature representation. An encoder-decoder based approach to alleviate the occlusion challenge was proposed in [62]. Region encoder was responsible for building a correlation among occluded and non-occluded regions. while the region decoder uses the spatial correlation to recover the occluded regions in particular frames. Moreover, in order to refine the spatial region feature completion a temporal region feature completion module was plugged in that was responsible to create the long-term temporal contexts.

Another network was presented in [63]. They have combined the intrinsic relationship between the tasks of person re-id and semantic segmentation to alleviate occlusion. Proposed network consists of three branches *i.e.* semantic, local and global branch. Local branch obtains the part level features, global branch was responsible to obtain features robust to occlusion and the semantic branch generates the foreground-background mask of a person image that basically leverages the non-occluded regions of human body. These three branches were trained collectively to obtain the discriminative representation of pedestrian image. In [55] a matching framework was proposed that consists of two modules *i.e.* a local and a global. Local matching was based on patch-level while the global matching was part-based to provide the spatial information. Presented approach was tested modified existing dataset (explicitly included partial person images) and a new partial person re-id dataset.

5.1.2. Attention-based approaches:

Occlusion as an explicit challenge was handled in [64]. Attention mechanism was used to learn more about non-occluded parts instead of occluded parts. At the matching stage only the visible shared regions based on pose landmarks are compared to obtain the efficient results without any manual cropping. Human pose landmarks worked as guidance to learn only the non-occluded regions. They have generated new occluded dataset named Occluded-DukeMTMC from DukeMTMC to test the proposed framework. While testing, both gallery and probe dataset contains occluded images to make it consistent with real world scenario. Another robust framework to resolve the challenge of partial occlusion was proposed in [65]. Spatio-temporal information was used explicitly to recover the occluded frames based on the visible body parts. Neighbouring frames were used to recover the information about occluded parts that had resulted to obtain accurate temporal information. A temporal attention layer was introduced that accurately learned the missing information from adjacent layers efficiently.

To address the problem of occlusion in more crowded situations a novel deep network named PISNet is presented in [66]. It is comprised of Query-Guided Attention Block that enhances the feature learning process of the target image in the gallery under the direction of query. Because of the improved location accuracy and attention based distinctive feature learning improved results are obtained even in occluded regions. Yingquan wang *et al.* in [67] targeted the challenge of occlusion. In the proposed novel framework both short and long-term temporal information was learned using attention mechanism in a hierarchical manner. Learned spatio-temporal representation is then aggregated by an

aggregation block that strengthens the learned features and makes them more discriminative to produce useful results on benchmarks.

In [68] a non-parametric attention mechanism was adopted that takes the video pairs as input and provides the matching score. Attention mechanism helped in refining the intra and inter sequence representation of input videos and outputs self and collaborative feature representation of each video. A generalized pairwise similarity measurement was also presented that calculates the pairwise similarity measurement of feature representation. Finally the matching result was obtained using binary classifier. Guangyi Chen *et al.* in [69] videos are sliced into different spatial-temporal units. Purpose was to preserve the body-structure information. Quality scores of these spatio-temporal units was obtained using attention model.

In [70] introduced discriminative explicit pathways to learn unique temporal and spatial features using existing 3D convolution network architecture. Each proposed path is responsible to learn motion (dynamic features) and appearance (static features) information of specific person. Moreover local and global features are also learnt to resolve the occlusion and misalignment. These features together make the approach unique and effective using attention maps. Improved results are obtained when tested on video based benchmark datasets. Another approach is [71] to handle the occlusion challenge based on Spatial and Temporal Memory Networks (STMN). Spatial part stores the information about spatial distractors while temporal memory stores the attention information in person videos. Based on collected information aggregated features are learned to obtain the representation that itself is refined at frame-level of person.

5.1.3. Self-Attention based approaches:

An end-to-end part aware transformer was proposed in[72] in a weakly supervised manner. Self-attention mechanism was used to capture the context information of full image. Transformer based encoder decoder architecture was used for pixel based context and diverse part-level discovery.

Fig. 7 summarized the reported SOTA results on occlusion. While Table 5 gives a brief overview of the published work done on occlusion challenge.

Table 5: Results obtained on occlusion challenge against each dataset. Results in bold are the highest.

SN	Paper	Dataset	R1/mAP	Code availability
1	STRF, 2021, [70]	MARS DukeMTMC-VidReId iLIDs	90.3/86.1 97.4/96.4 89.3/–	No
2	STMN, 2021, [71]	MARS DukeMTMC-VidReId iLIDs	90.5/84.5 97.0/95.9 82.1/69.2	Yes
3	PSTA, 2021, [67]	MARS DukeMTMC-VidReId iLIDs PRID-2011	91.5 /85.8 98.3 /97.4 91.5/– 95.6/–	Yes
4	SGPR, 2021, [59]	Market-1501 DukeMTMC-ReId Occluded-ReId Occluded-Duke	96.1/89.3 91.1/81.3 78.5/72.9 69.0/57.2	No
5	PE-PGFA, 2021, [60]	Occluded-ReId Occluded-Duke Partial-iLIDs	81.0/71.0 62.2/46.3 80.7/85.7	No
6	FPR, 2019, [56]	Market-1501 DukeMTMC-ReId CUHK-03	95.42/86.58 88.64/78.42 76.08/72.31	No
7	PGFA, 2019, [64]	Market-1501 DukeMTMC-ReId Occluded-Duke	91.2/76.8 82.6/65.5 51.4/37.3	No

8	AMC-SWM, 2015, [55]	Partial-ReId	53.14/–	No
9	PISNet, 2020, [66]	Market-1501 DukeMTMC-ReId	95.6/87.1 88.8/78.7	No
10	GASM, 2020, [57]	Market-1501 DukeMTMC-ReId MSMT-17	95.3/84.7 88.3/74.4 79.5/52.5	Yes
11	PAT, 2021, [72]	Market-1501 DukeMTMC-ReId Occluded-ReId Partial-iLIDs	95.4/88.0 64.5/53.6 81.6/72.1 88.0/76.5	No
12	HOReID, 2020, [58]	Market-1501 DukeMTMC-ReId Occluded-ReId Occluded-Duke Partial-iLIDs	94.2/84.9 86.9/75.6 80.3/70.2 55.1/43.8 72.6/85.3	Yes
13	VRSTC, 2019, [65]	MARS DukeMTMC-VidReId iLIDs	88.5/82.3 95.0/93.5 83.4/–	No
14	ATNet, 2019, [73]	Market-1501 DukeMTMC-ReId	45.1/24.9 55.7/25.6	No
15	IGOAS, 2021, [61]	Market-1501 DukeMTMC-ReId	93.4/84.1 86.9/75.1	No
16	SCAN, 2019, [68]	MARS iLIDs	87.2/77.2 88.0/89.9	No
17	STAL, 2019, [69]	MARS iLIDs	71.5/50.8 76.7/–	No
18	SRFC, 2021, [62]	MARS DukeMTMC-ReId CUHK-03 MSMT-17 MARS DukeMTMC-VidReId	95.2/89.2 90.7/80.7 78.0/81.1 82.0/60.2 90.7/86.3 97.6/97.0	No
19	SORN, 2021, [63]	Market-1501 DukeMTMC-ReId	94.8/84.5 86.9/74.1	No

5.2. Pose Variance

Despite of hundreds of research papers, person re-id is still challenging to be solved mainly due to complex view variations and pose variations in the person images as shown in Fig. 1. The problem of body misalignment caused due to pose/viewpoint variations and occlusion results in imperfect detection.

5.2.1. CNN-based Approaches:

Spatial Interaction and Aggregation (SIA) module was introduced in [74] to increase the feature learning capability of CNN. The module adapts receptive field according to pose and scale of input image and hence resolves the issue of large body variations. Another module named Channel Interaction and Aggregation (CIA) was also used to semantically aggregate the similar channel features to make the better feature representation even for small visual cues. In [75] part based model bottom-up approach was used in a fully conventional way to improve the long-range predictions. Once the localization of keypoints is done then greedy decoding process was used to group the instances. Proposed approach was confident in detection as the detection process starts from distinguishing key-points *i.e.* nose to produce best results even in clutter and with variant poses.

In [76] presented a novel Pose-driven Deep Convolutional model that learns the global and local representation simultaneously. Softmax loss was used to learn the global (whole body) representation while a Feature Embedding

subnet was implemented to learn the local (body-part) representation. This local representation makes possible to learn the affine transformation and relocate the regions for easy recognition across multiple cameras. Pose Transfer Network further makes it possible to handle the pose variations. Similarity measure is then fed with this effective fusion of features, to obtain the effective results. It is noted that model is end-to-end and model weights and representations are learned jointly.

Muhammed Kocabas *et al.* [77] has presented his work to learn multi-person pose estimation using 4x faster bottom-up approach. Their proposed model is multi-task. Person detection, segmentation and pose estimation are jointly learned using shared backbone. Network have the key-points and after person detection it forms the pose by assigning key-points information to each detected person. Their method can also be extended for other tasks *i.e.* person segmentation. Yeong-Jun Cho *et al.* in [78] analysed camera viewpoints and person poses. Captured images are calibrated and their respective pose is estimated. Proposed multi-pose model gives representative features that are clustered into four groups according to poses *i.e.* left, right, front and back. In a weighted summation aspect, matching scores are calculated between multi-pose models to produce effective results.

In [79] challenge of extreme changes in pose and view point is observed. A novel unified framework that combines the saliency and semantic parsing to enhance the performance results. Saliency is about focusing at the points that are at the first glance by human eye. Global representation is carried out by saliency and semantic parsing masks. These generate the two types of complementary feature maps that improve the results. Semantic parsing is used to encode all parts of the person and this results in overcoming the challenges of occlusion and misalignment in the bounding box. Hence each sub-net stream learns to resolve different scenarios of the problem. In [80] attribute based approach was used to approach the challenge of large visual variations and spatial shifts. Large spatial shifts are caused by different pose variations and camera views. Pedestrian attribute assisted CNN-based framework was consists of two parts. In the first part attributes are learned. LOMO features are extracted that are specifically in use to make the model viewpoint invariant. These hand crafted low level features are combined with high-level learned CNN features to obtain robust and diverse feature representation model. The learned embeddings from this first step was then used at second step. In this second step two neural networks fuse together. One is pre-trained network on attribute labels and the other is CNN pre-trained on person re-id labels. These networks are integrated into triplet architecture. Hence optimal fusion parameters are learned in this manner.

Niall McLaughlin *et al.* in [81] proposed a strategy to tackle the challenge of appearance changes. In the proposed unified framework invariant feature extraction and supervised learning are combined. To overcome the problem of over-fitting several techniques were used including multitask learning, data augmentation and dropout. Siamese network architecture was used to extract the useful features. Siamese architecture helps to train the network to generate low dimensional feature representation. Diverse images of same person are mapped onto same location in feature space whereas images of different people are mapped onto different locations in the feature space. This network setting helps to learn the subtle cues from diverse set of images. In [82] Jin Wang *et al.* proposed to learn the deep representation with a novel adaptive margin list-wise loss. In training ranking lists are introduced to replace image pairs. Ranking lists resolve the issue of data imbalance. And Adaptive margin in the list-wise loss function helps in assigning larger margin to hard negative samples. Four convolutional neural networks were combined in one architecture, each network takes input of different body parts or scales. In training stage network was jointly optimized with similarity layer that can be separated at testing stage; this separation accelerates the computation.

To overcome the challenge of pose variation ensemble of invariant features were proposed in [83]. Ensemble of invariant features helped to achieve the robustness against multiple challenges includes partial occlusion, camera color changes, pose and viewpoint variations. Ensemble of invariant features includes pose-invariant features of specific regions and of holistic image. Invariant features are extracted using DCNN. Region based features are extracted from different body parts via Gaussian Mixture Models on color histograms. Each Gaussian distribution represents the dominant color mode. Main goal of [84] was to come up with a solution that is robust in handling both pose variations and misalignment. Pose Guided Representation composed of two components *i.e.* Pose Invariant Features (PIF) and Local Descriptive Features (LDF). PIF handles the pose variations, it approximates pose invariant representation via pose estimation and normalization. While LDF resolves the misalignment errors, it predicts the discriminative representation learning via body region segmentation. To achieve the robust results at test time, the pose estimation and region segmentation only applies during training time.

Shoubiao Tan *et al.* addressed the challenge of variation in pose and viewpoint across multiple camera views in paper [85]. To model the variations of poses a ranking based strategy was proposed that fuses dense invariant features. Images are first divided into dense sampled patched inspired by the fact that local features perform better than global

features. To achieve viewpoint invariant representation, it is considered that discriminative parts of a person appear in different regions and in different views. Hence corresponding patch of first image is searched in the neighbourhood of second image. Moreover based on the dense invariant features Support Vector Ranking was used further to learn the transformation across the views and hence results are improved. A network based on Inception architecture was presented in [86] to handle multiple challenges *i.e.* pose, view and illumination. They have made the network to learn the structural aspects of human body. The network was trained and tested not only on real datasets but on a newly formed synthetic datasets as well.

In [87] teacher-student model was used in which teacher acts a guider to extract the global features. These global features are then used to train the other branch of the model that is local feature branch. Moreover a part prediction alignment module was also introduced that align the different parts of same person that resulted in effective estimation of person pose and alignment. To alleviate the pose variation problem a fine-grained person re-id solution was proposed by Jiahang Yin *et al.* in [88]. The proposed solution was focused on dynamic pose features. Two types of pose features are incorporated *i.e.* motion-attentive local dynamic pose feature and joint-specific local dynamic pose feature. Motion and global person features are incorporated in the model to obtain the discriminative features. Moreover a new dataset was also proposed to better resolve the issues specific to fine-grained person re-id problem.

In [89] a channel parse block was responsible to extract the pose information at pixel level. The required information was obtained by suppressing the background to avoid the inaccurate detection due to occlusion. Moreover a patch level alignment mode was also proposed to handle the misalignment at local level in a fine-grained way. In [90] a high order re-id framework was proposed to handle the pose alignment problem using semantic fine-grained part details of multi-level feature maps. These multi-level feature similarity maps were then used to find the difference of similarity among aligned and misaligned parts of person images. As the similarity information among two person images found reduced so the proposed framework successfully increase the robustness of pose features to rectify the misaligned pose information.

A unified end-to-end trainable framework to handle the challenge of pose and viewpoint variation was proposed in [91]. For accurate ranking of gallery images a novel Kronecker Product Matching operation and group shuffling random walk operation was presented. Together these modules are responsible to accurately rank the probe to gallery and gallery to gallery affinities. This accurate ranking resulted in improved learning of features when embedded with deep learning frameworks. Mid-level multi-type human attributes were learned in [92] in weakly supervised manner. These human attributes were then used to learn features of images with large visual variance for person matching. Contextual cues among attributes were considered to play the vital role in boosting the accuracy of the model in a progressive way.

In [93] (both for image and video) challenge of pose variation and alignment was targeted using both fine and coarse information of acquired person image. They have adopted an unsupervised re-ranking framework, produced results are re-ranked based on simple aggregate of Euclidean distances among the provided gallery and probe images. Pose Invariant and Embedding (PIE) [94] resolves the challenge of pose changes and errors in pose estimation. At first place PosBox is designed to formulate a bounding box around each individual. Due to the careful use of pose estimators [95], PoseBox produces well-aligned person images so that discriminative and up to the mark features could be learned even in intensive scenarios of pose changes. Secondly Pose Box Fusion (PBF) CNN is added during formulation of PoseBox to reduce the impact of information loss and errors in pose estimation. PBF-CNN takes three streams of input *i.e.* Original image, PoseBox and confidence score of pose estimation. PBF achieves globally optimized results on PoseBox as compared to original image. This PIE plays its role as FC activation of PBF network.

Pose variation is a complex challenge and is targeted in [96], GAN framework for transferable pose variation was based on sample augmentations. Xuecheng Nie *et al.* [97] challenged pose variations to achieve efficient results on person re-id. They have tried to overcome the limitations of top-down and bottom-up approaches by using unified model based on regression process. The proposed joint framework simultaneously models the detection and joint partition. This has helped to precisely detect with partitioned joints using candidate voting scheme, that further speed-ups the pose estimation. Feed-forward pass helped in achieving higher results as compared to other top-down approaches for achieving optimal improved performance results on multi-person joint configurations pose estimations. Another representation learning scheme was proposed in [98] for video person re-id that is based on adaptive graphs. Pose alignment connection was built to have an adaptive structure aware graph. In order to refine the regional features, feature propagation was performed on adjacency graph in an iterative manner. To make the representation more discriminative temporal resolution aware regularization scheme was also evolved that ensures the consistent temporal resolution for alike identities.

Xiaoqiang Hu *et al.* in [99]. In the presented hypergraph video pedestrian re-id method a posture structure and action relationship was explored. It makes full use of the walking posture of a pedestrian. They have used graph convolution network to preserve the structure information obtained from a pedestrian image. Structural relationship was then formed by detecting the joint point regions of the pedestrian. Moreover action hypergraph was also constructed by using the action information obtained by tracking the moving position of the detected joint points. Structure and action information together forms the saliency score which is then converted into a distribution problem to obtain the final results. In [100] a unified procedure is describe to handle the challenge of pose variation by stripping down image horizontally in spatial domain to maintain the geometric structure. Feature maps are obtained using [101]. Feature maps of each stripped part were compared with corresponding part in other image on the basis of learnt similarity measure. Each similarity measured score is combined to achieve flexibility. This approach is more robust to occlusion and hence less risky of mismatching.

Le An *et al.* [102] challenged the pose variations caused due to multiple reasons *i.e.* illumination changes, appearance changes and occlusion. For re-id gallery and probe images are projected onto RCCA subspace. In RCCA subspace reference descriptors of probe and gallery images are generated via measuring the similarity among images and the reference data. Identity of probe image is determined via comparison of probe reference descriptor and gallery reference descriptor. Saliency based matching is further used to add re-ranking step that improve the results further. Yifan Sun *i.e.* in [103] designed a discriminative feature descriptor named as Part-based Convolutional Baseline that is declared general as it is able to handle the challenges like pose estimation, human parsing and uniform part partitioning. Other component named as refined part pooling is based on feature descriptor and is used to precisely locate the parts. Their idea is based on within-part consistency which states that pixels in a well located part are more alike to each other and more dissimilar from other pixels of other parts. And if a pixel in a part is more alike to the pixels of other parts then it is considered as outlier and this refers to inefficient partitioning. Refined Part Pooling (RPP) handles this outlier issue by reallocating the pixel to its closest matching part. RPP trained in semi-supervised way and needs no part labeling. This approach increased the generality of model and boosted the performance as well.

A pose-aware multi-shot matching strategy was presented in [104]. It estimates the pose information using pose estimation methods and then multi-shot matching is performed. The approach analysis both camera viewpoint and person poses. Four image clusters are formed on the basis of front, back, left and right poses. Multi-pose model is generated based on the four feature descriptors formed on the basis of image clusters. Matching scores are then calculated between generated multi-pose models. Additional cues like person poses and 3D scene information makes the proposed model more tractable.

5.2.2. Attention-based Approaches:

Challenge of pose variation was addressed [105] in an end-to-end way using overlapped activation penalty. Activation penalty pays more attention to the less activated regions of the image. Diverse complementary features were learnt effectively with less penalty. Challenges of background clutter, alignment, occlusion and pose are resolved in a unified framework [106]. In the proposed framework attention-aware pose guided network is used to handle the challenge of occlusion, misalignment and pose. And noise was excluded via filtering finer parts based on attention mechanism. Attention mechanism was adopted in [107] to tackle the problem of pose variation via efficiently capturing the global structure information to better formulate the attention based discriminative feature learning for person re-id task. Global structural information that is kind of clustering information is better to have the contextual information. This global relation helps to leverage discriminative attention on human body regions. In this way both global (semantic) and local (human body regions) are learned to produce the effective results.

An end-to-end Comparative Attention Network (CAN) based framework presented in [108]. Model is inherent with comparative components. They have tackled appearance variances caused due multiple factors mainly includes different poses and camera views. Proposed model adaptively finds the multiple local regions comprised of discriminative information via set of glimpses. At each glimpse (or region of attention) model creates different parts. Model takes location of previous glimpses and raw person images as input, and formulates next glimpse of local regions as output. These glimpses formulate a kind of dynamical pooling feature, resulted in enhanced performance as compared to conventional pooling features. Which is then integrated to further improve the performance.

An end-to-end pose-guided framework was introduced in [109] to tackle the challenge of occlusion. In the jointly learned visibility prediction model pose-guided attention and part visibility models are combined. Graph matching technique was then adopted to match the self-learned features of visibility cues and self-mine visibility score of the gallery and test image accordingly. Matching holistic images with occluded images is an inefficient approach and leads

to wrong results as well. Hence occlusion was dealt explicitly. Fig 8 shows the progress on pose variation in literature since 2015.

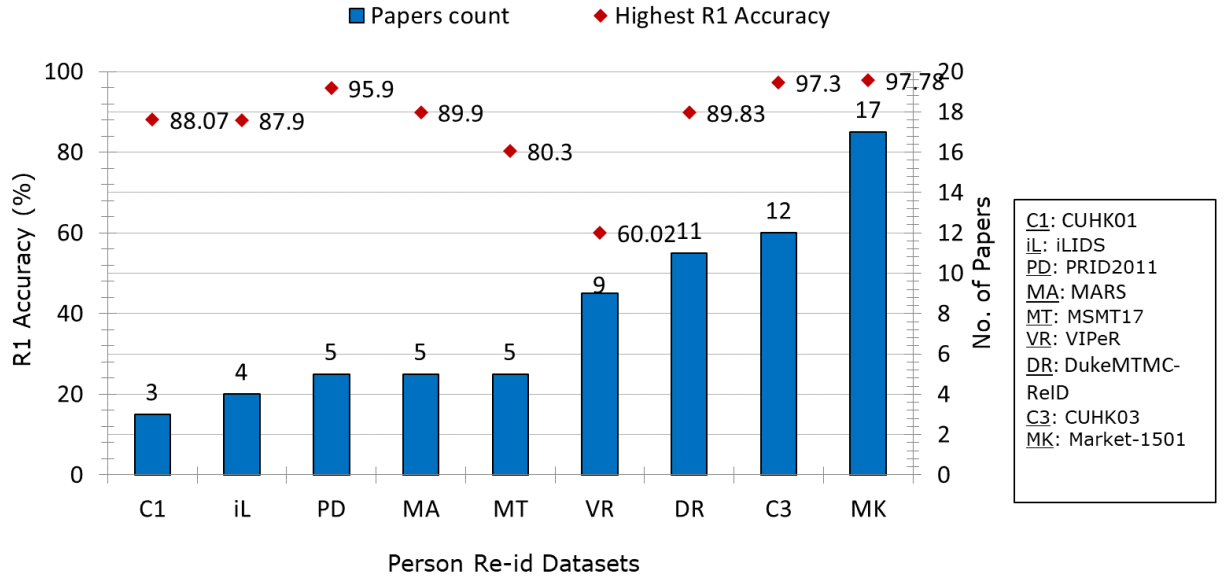


Figure 8: Progress on the challenge of pose variation for person re-id benchmarks

Table 6: Results obtained on Pose-variation challenge against each dataset. Results in bold are the highest.

SN	Paper	Dataset	R1/mAP	Code availability
1	PDC, 2017, [76]	Market-1501	84.14/63.41	No
		CUHK-03 (Labeled)	88.7/–	
		CUHK-03 (Detected)	78.29/–	
		VIPeR	51.27/–	
2	RGA, 2020, [107]	Market-1501	96.1/88.4	Yes
		CUHK-03 (Labeled)	81.1/77.4	
		CUHK-03 (Detected)	79.6/74.5	
		MSMT-17	80.3/57.5	
3	SIA+CIA, 2019, [74]	Market-1501	94.4/83.1	No
		DukeMTMC-ReID	87.1/73.4	
		CUHK-03 (Labeled)	92.4/–	
		CUHK-03 (Detected)	90.1/–	
		MSMT-17	75.5/46.8	
4	CAM+RAM, 2019, [105]	Market-1501	94.7/84.5	No
		DukeMTMC-ReID	85.5/72.9	
		CUHK-03 (labeled)	70.1/66.5	
		CUHK-03 (Detected)	66.6/64.2	
5	AACN, 2018, [106]	Market-1501	88.69/82.96	No
		DukeMTMC-ReID	76.84/59.25	
		CUHK-01(Detected)	88.07/–	
		CUHK-03(Labeled)	81.86/81.61	
6	PSE, 2018, [93]	Market-1501	90.3/84	Yes

		DukeMTMC-ReId	85.2/79.8	
		MARS	76.7/71.8	
7	PaMM, 2016, [78]	iLIDS PRID-2011	30.3/– 45.0/–	No
8	SCSR, 2016, [100]	ViPeR Market-1501	53.54/– 51.9/–	Yes
9	SOMASet, 2018, [86]	Market-1501 CUHK-03 (Labeled)	81.29/56.98 85.9/–	No
10	PPA, 2021, [87]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	92.4/79.6 85.1/71.8 69.2/66.3 65.5/62.4	Yes
11	PIE, 2019, [94]	Market-1501 DukeMTMC-ReId CUHK-03(Detected)	87.33/69.25 80.84/64.09 45.88/41.21	No
12	CAN, 2017, [108]	Market-1501 CUHK-01(Labeled) CUHK-03(Labeled) CUHK-03(Detected) ViPeR	72.1/47.9 87.2/– 77.6/– 69.2/– 54.1/–	No
13	PaMM, 2018, [104]	MARS iLIDS PRID-2011	66.3/– 57.3/– 79.4/–	No
14	PAFAM, 2020, [98]	MARS iLIDS PRID-2011	89.8/81.1 84.5/– 94.6/–	Yes
15	FGSAM, 2020, [89]	Market-1501 DukeMTMC-ReId	91.5/85.4 85.9/74.1	No
16	HOReID, 2021, [90]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) MSMT-17	97.78/93.94 89.83/82.16 96.12/– 78.42/54.77	No
17	PBCNN, 2018, [80]	CUHK-03(Detected)	65.0/–	No
18	PCB, 2019, [103]	Market-1501 DukeMTMC-ReId CUHK-03 MSMT-17	93.8/81.6 84.5/71.5 63.7/57.5 69.8/43.6	No
19	KPMM, 2021, [91]	Market-1501 DukeMTMC-ReId CUHK-03	93.1/84.7 71.3/84.9 97.3 /96.4	Yes
20	PGR, 2019, [84]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected) ViPeR MSMT-17	93.87/77.21 83.63/65.98 92.15/– 89.61/– 60.02/– 66.02/37.87	No
21	PGR, 2019, [84]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected) ViPeR MSMT-17	93.87/77.21 83.63/65.98 92.15/– 89.61/– 60.02/– 66.02/37.87	No
22	FGPR, 2020, [88]	MARS	82.9/72.7	No

23	PA-HVPreid, 2021, [99]	MARS	89.9/79.6	No
		iLIDS	87.9/–	
		PRID-2011	95.9/–	
24	WSMTAL, 2017, [92]	Market-1501	56.6/31.2	No
		VIPeR	39.7/–	
		PRID-2011	24.2/–	
25	WSMTAL, 2017, [92]	Market-1501	56.6/31.2	No
		VIPeR	39.7/–	
26	PReID-RS, 2015, [102]	CUHK-01	31.1/–	No
		VIPeR	33.29/–	
27	DIFs, 2016, [85]	CUHK-01	39.46/–	No
		VIPeR	29.35/–	
28	SNML, 2016, [81]	VIPeR	33.6/–	No
29	DL-ReID, 2016, [82]	CUHK-01(Detected)	57.02/–	Yes
		CUHK-03(Labeled)	55.89/–	
		CUHK-03(Detected)	50.67/–	
		VIPeR	40.51/–	

5.3. Background Clutter

It's very basic and important problem in person re-id. Example of background variations can be seen in Fig. 1.

Complex background makes the detection of pedestrian difficult. Existing methods are not capable to comprehensively address this challenge mainly due to two reasons. Firstly, available datasets have multiple images of one person with same background taken by less number of cameras like dataset CUHK01 and CUHK03 were collected by using 2 cameras, and Market-1501 dataset collected by using 6 cameras. When deep learning based models trained on single datasets they perform usually poor on the other available datasets due to different backgrounds. Secondly, (in the cases where background is complex) most of the existing method ignore the background they just only focus on the visual appearance of the person. Progress of last 6 years is shown in Fig. 9.

5.3.1. CNN-based Approaches:

The challenge of variations in background was investigated in [110]. In the proposed framework human parsing maps are calculated to cater the background influence. New dataset with different backgrounds was created. Data augmentation technique was used on existing datasets to generate images. In [111] fully annotated person re-id is proposed that is divided into two phases *i.e.* detection and person re-id. In the first human ROIs are automatically extracted using adaptive Gaussian Mixture Model method and HOG-SVM detector. Casting shadows are removed using density based score matching. In this scheme both chromatic and physical based features of shadow regions are taken into account. Good performance results are obtained even in the presence of clustered background with occlusion.

In Deep Person model[112] contextual information of body parts *i.e.* from head to foot was integrated to strongly relate the one part with another. LSTM was used in an end-to-end manner, this resulted in enhancing the discriminative features of local parts that will align better to full person (due to prior information of body structure). Global full body representation was also adopted. For identification task both global and part-based features are fed into two separate network branches. Difference of combining features from other approaches lies in ranking task branch that used triplet loss to learn the similarity measure explicitly.

In [113] long-short temporal spatial clues are used to obtain the robust representation of features. Proposed network does so by combining the motion appearance and motion refinement features. Motion appearance provide the temporal clues (person specific features) by suppressing the background clutter present in multiple scenes. While motion refinement part (activates the person specific features) incorporates the motion refinement layers that are to be executed in multiple motion-excitation blocks of CNN. In this manner the network can differentiate the persons in multiple scenes. A self-paced learning framework was presented in [114]. In order to learn the discriminative features and distance metric, a self-paced constraint and a regularization technique was implemented. A part-based

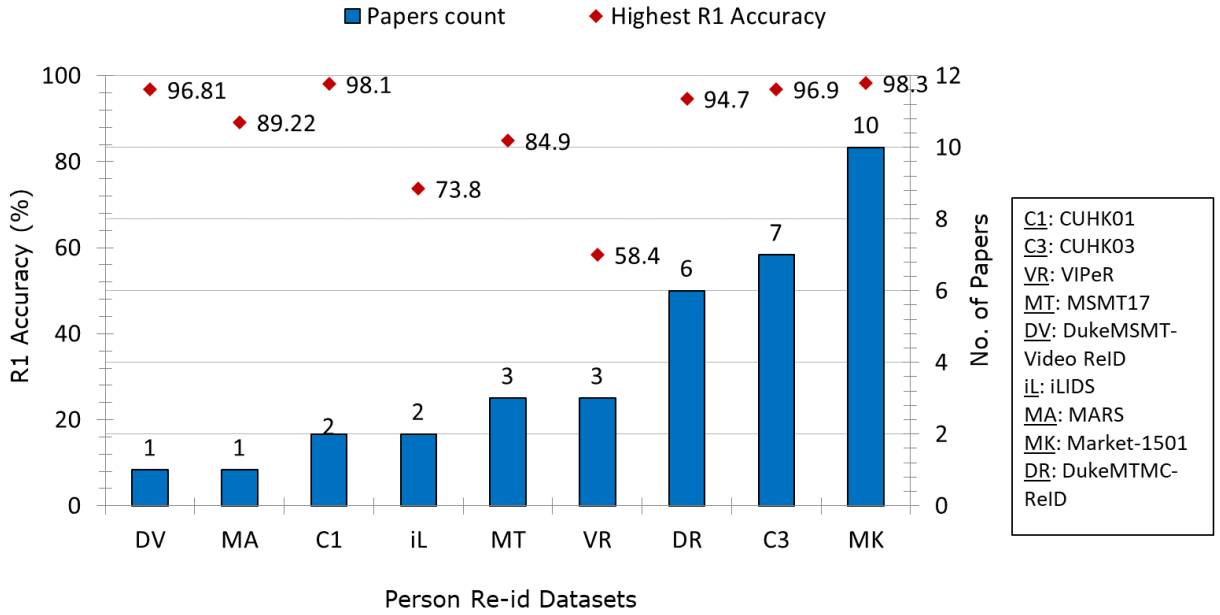


Figure 9: Progress on the challenge of background clutter for person re-id benchmarks

deep network was also built part features are learnt in lower convolutional layers and then embedded into higher layers to have the rich and discriminative features.

5.3.2. Attention-based Approaches:

In [115] an end-to-end Foreground Attentive Neural Network (FANN) is built to distinguish the individuals across non-overlapping camera views. In this attention based approach network focuses on foreground persons. Each image was first passed through encoder, decoder network. Encoder network extract features out of whole image, output of encoder is then further taken for learning discriminative features. Decoder network rebuilds the binary mask for each person in present in foreground of an image. Due to the regularization of decoder and use of local regression loss function encoder slowly start paying attention. Novel triplet loss is used to effectively learn the discriminative features. Use of triplet loss intra-class distances were minimized and inter-class distances were maximised at the same time.

An extended version of [115] was presented in [116]. It is an end-to-end foreground attentive neural network with symmetric triplet loss function. Framework is comprised of three sub-networks *i.e.* Foreground attentive sub-network, Body-part sub-network and Feature fusion sub-network. Foreground attentive sub-network comprised of encoder decoder network that takes input of images and focuses its attention on foreground part of input images. Body-part network takes encodes feature maps, slice them and learn features. Resulting feature maps are then fused in third and last sub-network. Finally normalized feature vectors are passed to symmetric triplet loss function. A soft mask based end-to-end foreground aware network was presented in [117]. In order to model the background both pedestrian and camera ID were used. A target attention loss helped in focusing on foreground pedestrian features by reducing the negative impact of changes in background. It is noted that as compared to existing approaches no additional dataset was required to train the model. Improved results have demonstrated the effectiveness of the proposed model.

A multi-level attention and fusion model was proposed in [118]. Multi-level attention module has helped in learning the global level features while the multi-layer fusion module has helped in increasing the feature expression at fine granular level. Implementation of the presented model has improved the results as compared to existing. To remove the background interference Xin Ning *et al.* designed a feature refinement approach in [119]. Instead of directly focusing on high response features, complete features of a person was extracted and then highly valuable features were identified using multi-branch attention network that has resulted in increased performance of the model on benchmark re-id datasets.

Table 7: Results obtained on Background-clutter challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	CAR, 2019, [115]	Market-1501 DukeMTMC-ReID CUHK-03 (Labeled) CUHK-03 (Detected)	96.1/84.7 86.3/73.1 96.9/– 93.2/–	No
2	PRGP-DNN, 2018, [110]	Market-1501 CUHK-01 (Detected) CUHK-03 (Labeled) VIPeR	81.2/– 80.2/– 92.5/– 51.9/–	No
3	FANN, 2019, [116]	Market-1501 DukeMTMC-ReId CUHK-01 (Labeled) CUHK-01 (Detected) CUHK-03 (Labeled) CUHK-03 (Detected) VIPeR	94.4/82.5 85.2/70.2 98.1/– 81.2/– 92.3/– 70.2/– 58.4/–	No
4	TEM-ReID, 2021, [117]	Market-1501 DukeMTMC-ReId MSMT-17	95.0/84.6 88.7/77.0 76.8/51.0	No
5	HOG-SVM, 2017, [111]	iLIDS	73.8/–	No
6	SBSGAN, 2021, [120]	Market-1501 DukeMTMC-ReId	87.9/80 79.7/71.5	No
7	LSTS-NET, 2020, [113]	Market-1501 CUHK-03(Labeled) CUHK-03(Detected) MARS DukeMTMC-VideoReID iLIDS	95.8/93.0 78.23/72.3 70.11/67.9 89.22/83.12 96.81/93.91 60.92/–	No
8	LSTM-PDReID, 2020, [112]	Market-1501 DukeMTMC-ReID CUHK-03(Labeled) CUHK-03(Detecte)	94.48/85.09 80.9/64.8 91.5/– 89.4/–	Yes
9	DSPL, 2018, [114]	Market-1501 CUHK-01 CUHK-03 VIPeR	87.05/– 81.33/– 73.16/– 56.32/–	No
10	MEMF, 2021, [118]	Market-1501 MSMT-17	96.11/89.45 82.89/59.8	No
11	WFCB-ReID, 2021, [119]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected) MSMT-17	98.3/94.2 94.7 /90.3 88.6 /84.9 84.2/80.6 84.9/66.7	No

5.4. Misalignment

In person re-id, body misalignment (the problem in which body parts of person are spatially misaligned) is the essential challenge. The form of it is shown in Fig. 1.

Human detection becomes imperfect when the body parts were not perfectly aligned and it also makes the matching of person difficult between the probe and gallery image in person re-id. To build an efficient model which provides supervision in part alignment is still challenging task. Fig. 10 describe the progress on body misalignment in each of the image based datasets.

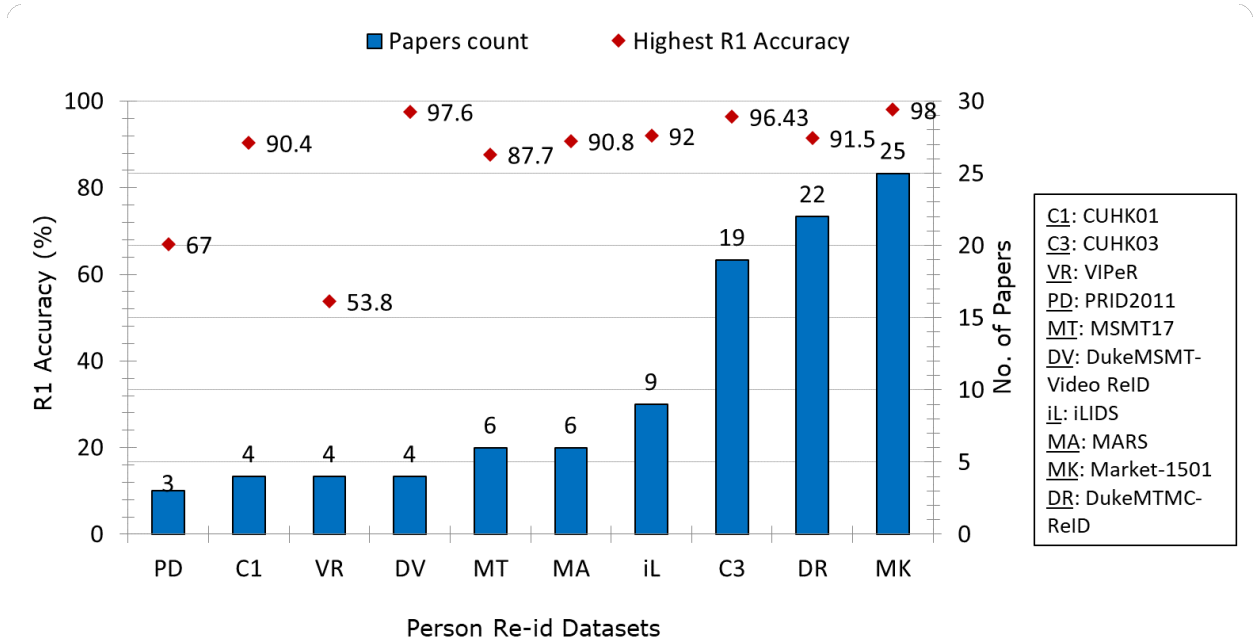


Figure 10: Progress on the challenge of misalignment for person re-id benchmarks

In past years, to resolve this challenge conventional neural network based approaches have been proposed and addressed it by using part based alignment and by focusing on shared regions in the images. As part based alignment has shown essential role in the generalization of re-id model.

5.4.1. CNN-based Approaches:

One of the framework to handle the challenge of misalignment was presented in [121]. Person image was divided into spatially semantically aligned 24 parts to learn the canonical surface based representation in UV space. Lingxiao He *et al.* [122] designed an end-to-end framework to leverage feature maps of different sizes without alignment problems as the matching is at pixel level. Yifan Sun *et al.* [54] proposed a supervised framework that consists of two sub networks with purpose of well aligned parts. One subnetwork was based on part-based convolution that generates a descriptor consists of uniform part-level convolutional discriminative features. Other adaptive subnetwork assigns outliers (out while doing uniform partitioning by first subnetwork) to nearest matching similar parts results in consistent refined parts.

In [123] part-maps are used to learn the human poses. These part-maps are then joined with appearance maps so that part-aligned representation can be computed. In the proposed model two streams are used to generate appearance maps and body part maps separately. Aggregation module then generates the feature maps that are part-aligned. As the computation of body-part information is not relative hence challenge of misalignment is reduced. In [124] structure of human body parts was used to better align the features of body region. These features are obtained using ROI pooling framework. And different semantic level features are obtained separately. These semantic level features are combined with local body region features using tree structured fusion network.

In [125] weakly supervised approach semantic parsing was used to address the misalignment. Proposed framework identifies the person body parts and belongings at pixel level to achieve aligned person re-id. Pseudo-labels of human body parts were generated and refined using iterative mechanism. The framework can effectively identify the occluded parts in an image using foreground and background clustering approach. In this way only features of visible body parts are learned that leads to fine feature matching. Yifan Sun *et al.* [126] proposed a fully convolutional network-ResNet-50 [127] named Visibility-aware Part Model (VPM) inspired from holistic person re-id [54] [128]. Proposed design was self-supervised, focuses the challenge of partial person re-id, misalignment and scalability. It learns the features of shared region among two images using region locator and region extractor. In VPM holistic images are used for end-to-end training while at test time distance/visibility score of shared region-to-region and whole image is computed to obtain the meaningful results on synthetic as well as realistic datasets.

In [129] Dynamically Matching Local Information (DMLI) aligns the horizontal strips without any explicit information of pose estimation. It also resolves the challenge of pose variation caused due to inaccurate detection, occlusion and change in viewpoints. Later they have combined the local branch that is DMLI with Aligned re-id++ (combination of local and global features) to learn the global features. Shortest path distance was used to align the local parts. In this manner local branch directs the global branch to learn more discriminative global features that results in improved performance. Irrespective of the fact that global level features have nearly global receptive field, they found that high-level feature maps are more suitable for aligning local parts.

A framework named RankSVM proposed in [130]. This patch matching based computational model handles the misalignment challenge caused due to factors like poses, viewpoints and illumination. This patch matching is integrated with saliency matching to increase the discriminative power and robustness. In this paper saliency refers to regions with distinct attributes and are reliable to find same person in multiple camera views. They have transformed the original high dimensional visual feature space into saliency feature space that is 80 times less in dimension, this helped to attain high efficiency and resolve overfitting. Features of each 10*10 patch were extracted using Dense color histogram and Dense SIFT that is computed around local region.

Sheng Li *et al.* in [131] addressed the misalignment caused due to pose variations and viewpoint. In the proposed cross-view dictionary learning model, representation power of learned features is improved by learning dictionaries at multi-level *i.e.* image-level, horizontal part-level, patch-level and is considered as a general solution to multi-level learning problem. Main idea is to exploit and captures the details at multiple levels that captures both local and global characteristic of an image, it jointly enhanced the performance results on challenges like misalignment, pose variations.

Changxing Ding *et al.* in [132] presented a multi-task part-aware network to handle the challenge of misalignment. Semantically aligned features of aligned parts are extracted in training phase. Model learns the representation specific to part by using regularization technique that result in selection of part specific channels. Global max pooling helped in making the learned features invariant to scale and translation. That fixed channels specified for each body part makes the approach robust at inference level to cope with blur and clutter may present in the image. Guanshuo Wang *et al.* in [133] focused on stripe-based approach to learn the features at multi-granular level. Instead of performing partition at input images or output features they have applied partition on intermediate representation to retain the local association among parts. Moreover, to deal with misalignment challenge random shifting augmentation was applied within bounding boxes appeared around detected persons in the image.

Poor spatial alignment for video sequences was targeted in [134]. Proposed framework consists of two modules *i.e.* temporal residual module and spatial-temporal transformer network module. Temporal module was responsible to obtain both generic and specific features present in consecutive frames while in second module semantic information of specific frame and its temporal context with adjacent frames were recorded. Model effectively align the person with major changes in appearance and hence outperforms the existing approaches. In [135] attributes are learned per patch. As a way local part features are extracted to handle the misalignment challenge. For final feature representation local features are fused with holistic features. hence effective results are obtained.

Wei Shi *et al.* in [136] a deep Image to Video re-id pipeline was proposed. Fine-grained features are learnt using a three dimensional semantic appearance alignment module. Module extracts the images that are aligned with local appearance. These aligned images were then aggregated with another multi-branch network that helped in weakening the influence of occluded body parts. Moreover, another module was responsible to handle the modality misalignment problem, it ensures the interaction among global representation of an image and video streams as a result discriminative fine-grained features were learned to produce the outperforming results. An unsupervised way to tackle the challenge of unaligned video based person re-id was presented in [137]. A video matching algorithm select and then match

the image sequences that are inaccurate or incomplete. Proposed approach does not rely on external labelling hence applicable to large scale unseen data to produce better results.

Patch wise metric learning approach proposed in [138]. Appearance measures of each patch was learned and then combined using deformable models. Patches were allowed to change the locations so to resolve the correspondence issue. Even if patch locations were found different but each have the same metric score hence the proposed method effectively helped to increase the training data and produce the effective results. In addition to spatial alignment temporal alignment was challenged in [139]. First they extracted the walking cycles (can be referred as gait) of a person in a chunk of video sequence. Each chunk was then divided into spatial and temporal data. Temporal sequence was then further divided into segments based on different walking cycle phases. While in the spatial domain image was divided into different body parts. Body-action unit was then formed on the basis of discriminative information obtained from temporal and spatial domains. Fisher vectors were extracted from body-action units to obtain a sort of generalized Bag-of-Words feature. Obtained features are then combined to form a vector that shows the visual/appearance of walking person.

Yang Shen *et al.* [140] presented a novel framework that handles the spatial misalignment caused due to multiple factors like camera-view and pose variations. Proposed framework learns the correspondence structure that include probabilities of matching patches among pair of cameras, this correspondence structure learning allows the better handling of misalignment. One-to-many graph in each correspondence structure of each patch allows to tackle the pose variations within each camera view, this graph shows the weights that depicts the matching probabilities of patches. Not only local but global context was also considered to achieve more reliable scores. Yang Shen *et al.* [140] presented a novel framework that handles the spatial misalignment caused due to multiple factors like camera-view and pose variations. Proposed framework learns the correspondence structure that include probabilities of matching patches among pair of cameras, this correspondence structure learning allows the better handling of misalignment. One-to-many graph in each correspondence structure of each patch allows to tackle the pose variations within each camera view, this graph shows the weights that depicts the matching probabilities of patches. Not only local but global context was also considered to achieve more reliable scores.

Inter and intra local relationship among extracted features was maintained simultaneously in [141] using part-guided graph convolution network. For optimization purpose two types of graphs were formed that describe the relationship among adjacent parts. First is inter-local graph of same parts of person image and other one is intra-local graph of variant parts of a person image. At the graph convolution operation was performed to inject the representation of person images. In [142] decoder and encoder based approach was used to handle the challenge of misalignment. Proposed model takes the advantage of both CNN and attention based architecture. On the basis of these architectures multi-grained spatio-temporal and positional spatio-temporal features are effectively learned to handle the misalignment.

5.4.2. Attention-based Approaches:

Wei Li *et al.* [50] proposed a framework to resolve the challenge of alignment and representation learning. Proposed framework jointly learns the hard region-level attention along with soft pixel-level attention in specified bounding boxes to efficiently learn the feature representation. In [47] misalignment and background clutter challenges are addressed. They have combined multiple modules into a single module to extract the better feature maps. Their proposed design is able to obtain both channel wise and position/spatial information using channel and position attention modules. In [143] a deep bilinear attention framework was formulated to handle the challenge of misalignment and representation learning. Two attention modules are introduced, decision for outer attention module *i.e.* where to focus is taken by inner attention module. The module uses channel wise second order information and hence interdependency among global and local features is formed, in the meanwhile spatial information is also preserved.

Jianyuan Guo *et al.* [144] presented a supervised human body parts parsing approach to achieve efficient results on a challenge of alignment in person re-id. In the proposed framework information from both human and non-human parts was extracted. Parsing model was used to extract human parts aligned information while self-attention mechanism was used to group alike pixels. Parsing and attention modules are then combined to learn discriminative features for accurate human parts and coarse non-human parts. A unified deep supervised network to resolve the challenge of misalignment was presented in [145]. Novel fully attention based block can be inserted into any convolution neural network to obtain the deep features. Channel-wise and spatial-wise attention information was extracted using fully attention block to better extract the useful multi-scale features. The problem of vanishing gradient was also resolved because of deep supervision. Body part misalignment problem was also handled in [146]. They have decomposed the body parts into

regions and computes the representation accordingly. This computation is discriminative in itself, hence helpful for person matching. The model learning was inspired by attention mechanism.

In [147] challenge of misalignment was addressed with the help of pose estimation in a multi-branch network architecture. In the proposed attention based network, representation of body-part and whole-body is learnt. This representation is then fused on the basis of their contribution towards feature matching. Intra-attention directs to precisely learn the discriminative features of whole-body and body-part images. Whereas intra-attention simultaneously learn the optimal feature representation and attention maps for whole body and interested part of body. In [148] a novel end-to-end model for part alignment problem is proposed. Model not only detects the particular body parts but it also extracts the discriminative representation at part-level. They have divided the image vertically and horizontally to obtain the structural information without any alignment and clutter issue. Vertical module detects the body-parts while horizontal module learns the part representation using attention mechanism.

In order to learn semantically aligned part-level features a simple batch-driven approach was proposed in [149]. Two modules were used *i.e.* a guided attention channel and pair of regularization term. First module highlights the channel responsible for each part of person image in the output of the deep network while the other regularizer term that maintains the consistency among batches and hence make the process coherent and robust. In [150] a novel triplet loss was proposed that not only considers the alignment but also pays attention to salient parts of the person image. It measure how much effort is deployed to align two distributions. Distribution of local parts were formed using attention technique. Weights are assigned to each part according to learned distribution. Novel loss term rectifies the assigned weights and hence an elegant solution to misalignment was presented.

An attention mechanism was adopted in [151] that exhibit its focus on body parts rather than background. Alignment of person images was learned using identification procedure. Proposed mechanism effectively located the person in an image by placing a bounding box based on attention mechanism. Xinqian Gu *et al.* in [152] tried to resolve the temporal appearance misalignment in video based person re-id. Their appearance preserving framework is comprised of two parts. First part preserves the appearance at pixel level and other one is 3D convolution kernel that helps to model the temporal information. To ensure the temporal appearance alignment adjacent feature maps are reconstructed according to cross-pixel semantic similarity. An attention mask was also learned that finds the unmatched regions among reconstructed and central feature map. Learned attention mask was then imposed to avoid error propagation.

Table 8: Results obtained on Misalignment challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	STRF, 2021, [70]	MARS	90.3/86.1	No
		DukeMTMC-VideoReID	97.4/96.4	
		iLIDS	89.3/–	
2	DenseIL, 2021, [142]	MARS	90.8/87.0	No
		DukeMTMC-VideoReID	97.6/97.1	
		iLIDS	92.0/–	
3	ABD-Net, 2019, [47]	Market-1501	95.6/82.28	Yes
		DukeMTMC-ReId	89.0/78.59	
		MSMT-17	82.3/60.8	
4	BAT-net, 2019, [143]	Market-1501	95.1/87.4	No
		DukeMTMC-ReId	87.7/77.3	
		CUHK-03(Labeled)	78.6/76.1	
		CUHK-03(Detected)	76.2/73.2	
		MSMT-17	79.5/56.8	
5	P2-Net, 2019, [144]	Market-1501	95.2/85.6	Yes
		DukeMTMC-ReId	86.5/73.1	
		CUHK-03(Labeled)	78.3/73.6	
		CUHK-03(Detected)	74.9/68.9	

6	PAHR, 2017, [146]	Market-1501 CUHK-03 VIPeR	81.0/63.4 85.4/90.9 48.7/–	No
7	BBA+PWM, 2015, [140]	iLIDS PRID-2011	44.3/– 64.1/–	No
8	AP3D, 2020, [152]	MARS DukeMTMC-VideoReID iLIDS	90.7/85.6 97.2/96.1 88.7/–	Yes
9	ISP, 2020, [125]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	95.3/88.6 89.6/80 76.5/74.1 75.2/71.4	Yes
10	PCB, 2018, [54]	Market-1501 DukeMTMC-ReId CUHK-03(Detected)	93.8/81.6 83.3/69.2 63.7/57.5	No
11	MANCS, 2018, [145]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	93.1/82.3 84.9/71.8 69.0/63.9 65.5/60.5	No
12	PABR, 2018, [123]	Market-1501 DukeMTMC-ReId CUHK-01(Labeled) CUHK-01(Detected) CUHK-03(Labeled) CUHK-03(Detected) MARS	95.4/93.1 84.4/69.3 80.7/– 90.4/– 91.5/– 88.0/– 84.7/75.9	Yes
13	VAPM, 2019, [126]	Market-1501 DukeMTMC-ReId	93.0/80.8 83.6/72.6	No
14	EANet, 2019, [153]	Market-1501 DukeMTMC-ReId CUHK-03(Detected)	94.6/85.6 87.5/74.6 72.5/66.8	Yes
15	DSA-reID, 2019, [121]	Market-1501 DukeMTMC-ReId CUHK-01 CUHK-03(Labeled) CUHK-03(Detected)	95.7/87.6 86.2/74.3 90.4/– 78.9/75.2 78.2/73.1	No
16	DSR, 2018, [122]	Market-1501	83.58/64.25	No
17	HA-CNN, 2018, [50]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	91.2/75.7 80.5/63.8 44.4/41 41.7/38.6	No
18	SPReID, 2018, [128]	Market-1501 DukeMTMC-ReId CUHK-03	94.63/90.96 88.96/84.99 96.22/–	Yes
19	Spindle Net, 2017, [124]	Market-1501 CUHK-01 CUHK-03 iLIDS PRID-2011	76.9/– 79.9/– 88.5/– 66.3/– 67/–	Yes

20	CDPM, 2020, [148]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	92.2/86.0 88.2/77.5 81.4/77.5 78.8/73.3	No
21	BCD-Net, 2020, [149]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected) MSMT-17	97/92.7 91.1/81.6 86.2/81.6 84.2/78.7 84.1/63.7	No
22	PCB, 2020, [150]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) MSMT-17	95.5/88.2 89.1/79.8 72.5/69.8 81.4/59.7	No
23	RMGL, 2020, [133]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled)	96.2/90.1 90.7/81.5 20.8/18.7	No
24	ST2N, 2019, [134]	MARS iLIDS	80.5/69.1 57.7/–	No
25	MPN, 2020, [132]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	96.3/89.4 91.5/82.0 85.0/81.1 83.4/79.1	No
26	ADP, 2018, [147]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	94.99/86.47 86.04/74.57 96.43/– 93.58/–	No
27	AlignedReID++ 2019, [129]	Market-1501 DukeMTMC-ReId CUHK-03(Detected) MSMT-17	92.8/89.4 86.2/82.8 67.9/70.7 69.8/43.7	Yes
28	APDR, 2020, [135]	Market-1501 DukeMTMC-ReId	94.4/90.0 87.3/83.2	No
29	3D-SAA+CMIL, 2021, [136]	MARS DukeMTMC-VideoReId iLIDS	81.3/72.6 82.8/81.0 54.7/–	No
30	PGCN, 2021, [141]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) MSMT-17	98.0/94.8 91.1/85.2 86.7/83.6 87.7/72.7	No
31	TS-DTW, 2017, [137]	iLIDS PRID-2011	31.5/– 41.7/–	No
32	DPML, 2017, [138]	CUHK-01 CUHK-03 VIPeR iLIDS	75.9/– 84.0/– 51.7/– 82.2/–	No
33	PAN, 2019, [151]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	82.81/63.35 71.59/51.51 36.86/35.03 36.29/34.0	No

5.5. Scale Difference

In real time scenario objects may appear in smaller or larger form. Their appearance may also depend on camera setting that makes the scale a complex challenge. Learning most discriminative features is the primary objective of re-id. These features should be computed at multiples scales so the re-id model become capable to differentiate two persons. Existing re-id model are based on fixed scale approach. A fixed scale representation makes most informative features blur due to this re-id model performance suffer. Existing approaches mostly resize all the pedestrian bounding box images into single scale. Fixed or single scale approach is not so optimal and explicitly multi scale representation becomes necessary as shown in Fig. 1.

People can be easily distinguished by using global features like gender and detecting local images patches but optimal matching only become possible if the features are computed at multiple scales and combined [154]. In open surveillance scene images are captures at an arbitrary scale, this make it challenging to learn correlations among features of different scales. [155] based on Siamese network and capable to learn most discriminative and informative features at different scales and evaluate their importance for cross-camera matching. In past five years progress achieved on this challenge is shown in Fig. 11

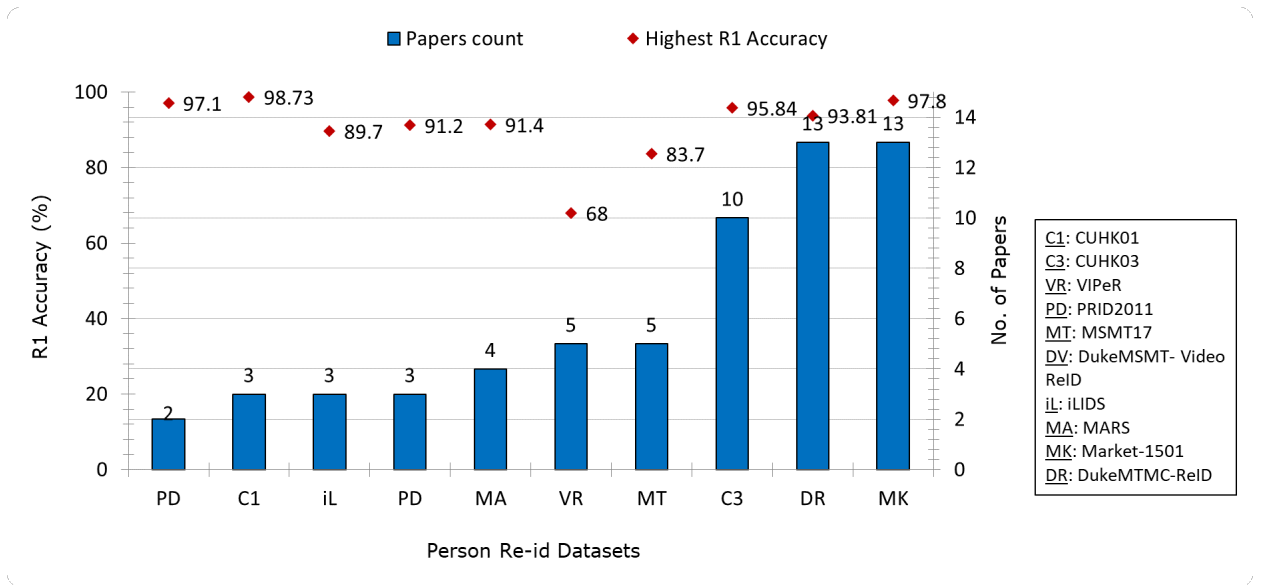


Figure 11: Progress on the challenge of scale variations for person re-id benchmarks

5.5.1. CNN-based Approaches:

A coarse-to-fine model [156] based on 3 dimensional sub-maps captures the discriminative information at different scales hence scalability challenge is addressed effectively. In [157] both inter-class and intra-class variations are taken as challenge and resolved using omni-scale feature extraction based on different receptive fields. Proposed network is lightweight due to usage of point-wise and depth-wise convolutions instead of standard convolution. Unified Aggregation (AG) Gate dynamically combines the generated channel wise weights of multi-scale feature maps. AG share parameters along all the streams and hence network enables to learn discriminative features at heterogeneous scales.

In [154] multi-scale feature learning problem is resolved using an end-to-end Deep Pyramidal Feature Learning CNN architecture. Model concurrently learns the scale-specific deep features. Effective results are obtained using multiple classification losses. There are scale specific branches that are correlated with each other (to maximise the scale specific discriminative features) and regularisation mechanisms that makes the model more effective for multi-scale issues. In [158] issue of multi-scale along with pose estimation was resolved using novel Spatial-Temporal Correlation and Topology Learning framework (CTL). The work was done for video-based person re-id. In order to obtain the diverse discriminative semantics local features were learnt at multi-granularity levels.

For video-based person re-id a novel 3D multi-scale convolution layer was introduced in [159] that can be inserted into any existing 2D convolution network. Depending on the location of insertion of 3D layer it splits into two variants *i.e.* local and global. Local learns the spatial-temporal cues among adjacent 2D feature maps while global learns the temporal relation among features of adjacent frames at global level. Together these features form the strong multi-scale combination and hence model got the immersive capability to learn discriminative features. An unsupervised generalized cross-dataset omni-scale approach was presented in [160]. In the proposed framework an omni-scale network learn the features at multiple spatial scales with assigned channel-wise weights to produce outperforming results.

In [161] a novel similarity learning model was presented. In the model they have combined the feature optimization using multi-view visual words and metric optimization. k-means clustering method was used to capture the multi-view visual words. A similarity function was then used to combine the common subspaces. Ancong Wu *et al.* [162] presented the scalable adaptive framework to address the challenge of scalability in unsupervised manner. In the presented approach source dataset was trained on labelled dataset however no labelling was done for practical testing. Teacher student transfer learning technique was adopted to obtain the effective results from trained model.

5.5.2. Attention-based Approaches:

In [163] visual similarities among images have been measured at different scales in an end-to-end manner. Moreover meaningful information have been extracted from feature map of an image using attention based spatial transformer Siamese network. A multi-scale model was presented in [164] that is cross attention-based. It is able to learn the discriminative information of different body parts of specific identity from multiple views. A new large scale Bird-View dataset was also proposed and tested along with existing benchmark datasets.

In [165] deep re-id network is proposed. Two novel layers are introduced to handle multi-scale challenge. Multi-scale deep learning layer learns the discriminative features at multiple scales. And leader based attention learning layer takes information of multiple scales and use this to selectively learn the optimal weightage assigned to each scale. Moreover pair of classification losses are used that strengthen the process of feature learning at multi-scale both at local and global level. Proposed methodology was useful due to the use of features that were extracted at multiple scales and locations and this makes the model generalised as some of the features become transferable.

Wei Zhang *et al.* in [166] proposed a framework for video based person re-id. In presented framework local regions are taken into account at multiple scales. Weights are assigned to local regions using attention mechanism at both spatial and temporal levels. These local features are aggregated to form rich representation of video, resulted in effective final outcome. A video-based person re-id approach was presented in [167]. A thorough spatial and temporal representation was obtained using two branches *i.e.* pyramid dilated convolution and pyramid attention pooling. Pyramid based strategy helped in extracting multi-scale features that has helped to mitigate the quality problems might present in the video *i.e.* partial occlusion. In [168] multi-scale pooled regions are fed into a novel deep architecture to extract the discriminative features at multiple scale semantic levels. Approach was made possible by using attention mechanism and was inspired by pyramid based methods.

Another multi-scale attention pyramid method to mitigate the scale challenge was presented in [169]. First features were divided into multiple local parts and then learned using attention mechanism. Then these features were merged and stacked using residual connection to form an attention pyramid. This attention pyramid was implemented in both channel-wise and spatial attention modules and hence reported the outperforming results as compared to existing state of the art methods. Another approach to learn complementary features was presented in [170]. Features were learnt from deep to shallow layers in a progressive manner. To focus on the layer specific features a two stage attention was also introduced to filter the noisy feature maps.

An end-to-end approach [171] to extract holistic and local feature maps using multi-scale omnibearing attention network. Multi-sized convolutions were used to obtain the local and holistic feature maps. Two kinds of attentions *i.e.* spatial and channel were also incorporated to extract the comprehensive feature representation.

Table 9: Results obtained on Scale challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	BV-Person, [164]	2021, Market-1501	96.0/89.2	No

		DukeMTMC-ReID	90.5/80.6	
2	OSNet, 2019, [157]	Market-1501 DukeMTMC-ReID CUHK-03 (Detected) VIPeR MSMT-17	94.8/84.9 88.6/73.5 72.3/67.8 68.0/– 78.7/52.9	Yes
3	MSDL, 2017, [155]	CUHK-01 CUHK-03 VIPeR iLIDS PRID-2011	79.01/– 75.64/– 43.03/– 41.0/– 65.0/–	Yes
4	DPFL, 2017, [154]	Market-1501 DukeMTMC-ReID CUHK-03(Labeled) CUHK-03(Detected)	92.3/80.7 79.2/60.6 86.7/82.8 82/78.1	No
5	CTL, 2021, [158]	MARS iLIDS	91.4/86.7 89.7/–	No
6	CFPModel, 2019, [156]	Market-1501 DukeMTMC-ReID CUHK-03(Labeled) CUHK-03(Detected)	95.7/88.2 89.0/79.0 78.9/76.9 78.9/74.8	No
7	AKA, 2019, [162]	Market-1501 DukeMTMC-ReID	49.7/24.6 47.6/31.1	No
8	STNs, 2018, [163]	CUHK-01 CUHK-03(Labeled) CUHK-03(Detected) VIPeR	88.2/– 87.5/– 86.45/– 50.1/–	No
9	MSTA, 2020, [166]	MARS iLIDS PRID-2011	82.28/69.42 70.1/– 91.2/–	No
10	DPRM, 2021, [167]	MARS DukeMTMC-VideoReID	89.0/83.0 97.1/95.6	No
11	PyrAttNet, 2020, [168]	Market-1501 DukeMTMC-ReID CUHK-03	97.8/95.8 93.0/90.9 86.8/88.0	Yes
12	M3D-CNN, 2020, [159]	MARS DukeMTMC-VideoReID	88.87/85.46 95.49/93.67	No
13	APNet, 2021, [169]	Market-1501 DukeMTMC-ReID CUHK-03 MSMT-17	96.2/90.5 90.4/81.5 87.4/85.3 83.7/63.5	Yes
14	PFE, 2021, [170]	Market-1501 DukeMTMC-ReID CUHK-03 MSMT-17	95.2/87.5 89.2/77.1 74.0/71.1 82.0/56.2	No
15	PREST, 2021, [172]	Market-1501 DukeMTMC-ReID	82.5/62.4 74.4/56.1	No
16	MuDeep, 2020, [165]	Market-1501 DukeMTMC-ReID CUHK-01	95.34/84.66 88.19/75.63 98.73/–	No

		CUHK-03 (Labeled)	95.84/–	
		CUHK-03 (Detected)	93.70/–	
17	OSNet, 2021, [160]	Market-1501	94.8/86.7	Yes
		DukeMTMC-ReId	88.7/76.6	
		CUHK-01	86.6/–	
		CUHK-03 (Labeled)	72.3/67.8	
		VIPeR	68.0/–	
		MSMT-17	79.1/55.1	
18	IALM, 2020, [161]	CUHK-01	68.44/–	No
		VIPeR	56.32/–	
19	MOAN, 2020, [171]	Market-1501	97.45/96.42	No
		DukeMTMC-ReId	93.81/92.82	
		CUHK-03 (Labeled)	90.07/90.32	
		MSMT-17	81.53/58.02	

5.6. Viewpoint Variance

Viewpoint is most important challenge in person re-id problem because different views of a pedestrian across non-overlapping cameras contain different information. Viewpoint support the learning algorithm specifically in identifying the pedestrian. Finding the particular angles are important for learning models to identify the person. To develop a model with outstanding generalizability is difficult because the visual appearance of the same person varies across different views across multiple cameras. Recent studies mainly focus on the learning of view invariant features. They initially extract view-generic features and after that view invariant model is learned. It is to reduce the distance between the intra class subjects and increase the gaps between the inter class pedestrians. Fig. 1 shows the viewpoint variation challenge.

But there are limitations in existing approaches, when the complex changes occur in the visual appearance between non-overlapping cameras, the view generic features become inadequate in solving person re-id task. Another problem in resolving this challenge is the lack of data and some datasets contain fixed and insufficient distribution of environmental factors *e.g.* in pedestrian viewpoint, some angles might contain few or zero samples. Due to importance of this challenge each year many papers were published in journals and conferences as well. The progress of work done so far on viewpoint challenge is shown in Fig. 12

5.6.1. CNN-based Approaches:

Xiaoxiao sun *et al.* [173] designed the practical solution by formulating the synthetic dataset PersonX for subjective study of various challenges *i.e.* occlusion, viewpoint variations, illumination changes, poses and various backgrounds. They have used IDE+ with ResNet-50 [127] as backbone with 36 angles and pre-trained ImageNet weights [40]. They have concluded with controlled experimentation that person with side views makes better query.

In [174] single dictionary was learned for both gallery and probe images simultaneously to obtain the view-invariant feature vectors of different people. To have discriminative dictionary encoding, model is then explicitly trained by applying constraint on association of sparse representations of the feature vector. At test time gallery images are compared (using Euclidean sense) to find the closest sparse match representation. In [175] efficient feature representation was learned. Presented method analyzes the presence of local features in horizontal direction. Occurrences of local features are maximized to make the approach viewpoint invariant. Further they have applied Retinex transform to reduce the impact of illumination and have reported the improved results.

In [176] visual variety was handled using teacher-student framework, teacher guides the student regarding multiple views and as a result student resulted in state-of-the-art in Image-To-Video by large margin. In paper [177] a novel cross-view semantic projection learning algorithm was proposed to model the feature transformation for person re-id. It retrieves the latent intrinsic invariant representation of persons. Three components are learnt simultaneously *i.e.* shared basis matrix, pair of semantic projection functions and optimal association function. In training phase the shared basis matrix explores the intrinsic structure of raw descriptors that are from different camera views. The projection function maps the original hand crafted features into common semantic space in the training phase. And the optimal association function captures the best association between semantic representation of alike individuals from multiple cross-views. Algorithm was then also generalised to multiple views datasets in the same paper.

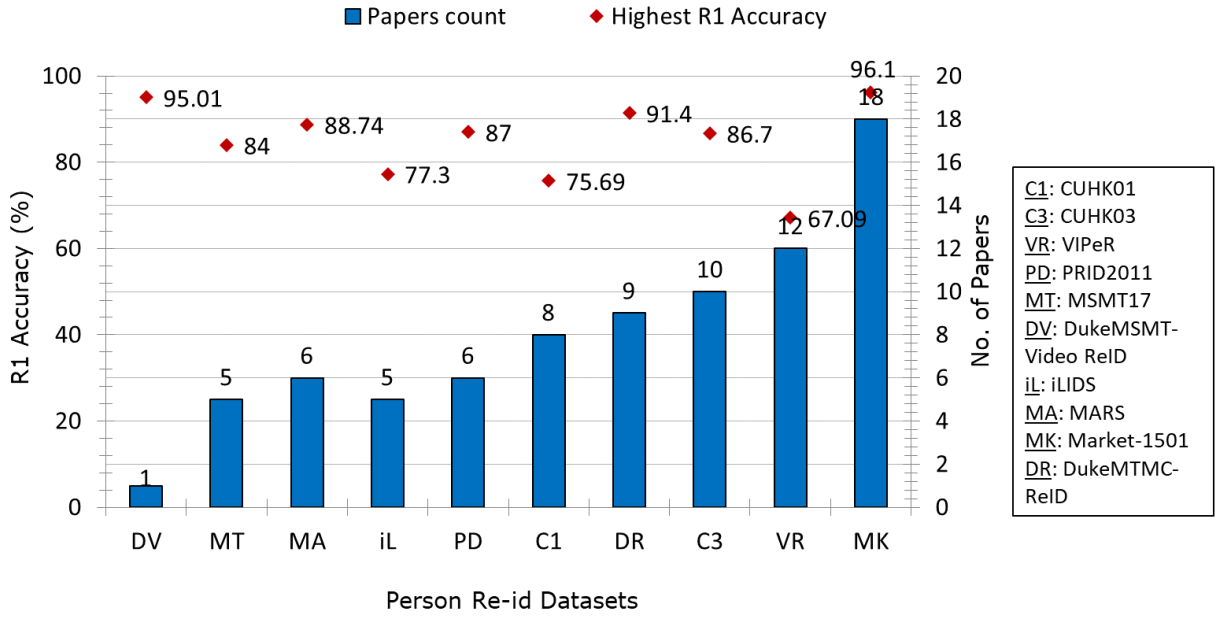


Figure 12: Progress on the challenge of viewpoint variation for person re-id benchmarks

Ying-Cong Chen *et al.* proposed a asymmetric distance learning model [178] to tackle the multi-camera view challenge. It transforms the different features that corresponds to different views to a common space. The model learns the camera-specific projection to resolve the problem feature discrepancy. Furthermore, consistency regularization helped in modeling the correlation among different views. Ziyan Wu et al. in [179] proposed a metric learning based viewpoint invariant algorithm to re-identify persons in cameras with disjoint field of view. It takes into account pose information prior from training data. After rectification color and texture histograms were used as descriptors to feature extractor from an image that was divided into six strips. Similarity measure is then used to assign the probe image to target class.

In [180] viewpoint specific approach named CRAFT was used for both view-generic and view-specific scenarios. Model was capable to adopt camera view features based on cross-view correlation and adaptive feature augmentation, this helps in transforming original features to new augmented space. Via this augmented framework view-generic can be induced to view-specific sub-models. To control the degree of correlation among sub-models camera view discrepancy regularization was also introduced in the CRAFT framework. To make the approach view invariant and generic to target datasets potential of deep learning techniques were explored and experimented. A view-specific model presented in [181] to handle the intra-class variations caused by viewpoints variations. Its two main contributions were: Cross-View Euclidean Constraint (CV-EC) and Cross-View Center Loss (CV-CL). CV-EC reduces the distance between features of same person from non-overlapping camera views. Whereas CV-CL was integrated to gain the high discriminative ability of view-specific deep networks. Proposed architecture can also deal with the applications using more than two cameras *i.e.* multi-view version.

One of the sever difficulty in re-id arises due to view variations in non-overlapping camera views. Hence a robust and discriminative descriptor is required to handle such challenge and it is resolved in [182]. This [182] novel deep learning to rank framework reduced the ranking cost of gallery images. In training phase, ranking model forms the relationship among input image pairs and their similarity scores via joint representational learning directly from raw image pixels. CNNs are used to form the relation between pair of images and their similarity score. In CNN feature representation and metric learning are integrated and generate the similarity score directly without relying on separate Euclidean or cosine distance measurement. Hence no need to compute feature representations separately. Network then learns the transformation in each ranking unit that tends to assign the highest score to the true match.

In [183] multiple metric learning method was presented that was based on cross-view quadratic discriminant analysis algorithm. The importance of each feature was decided using the discriminative power of each feature. Weights of all features are learned simultaneously using SVM learning framework. In this way gallery images were ranked based on the maximum attained score. Lin Wu *et al.* in [184] proposed a deep hashing framework that was fast and contains the discriminative representation of pedestrians. Hamming distance was then computed to rank the similar images closer to each other. Hash codes were assigned to each cluster of similar images to make the retrieval process faster. In [185] multiple convolution feature were extracted on the behalf of different body parts. Inter-layer interactions were also exploited to obtain the discriminative representation of person identities.

Alessandro Borgia *et al.* presented a metric learning approach in [186]. More discriminative feature space was learned using novel loss function. Inter-class and intra-class relationship of person identities was steered by meta center term and centers dispersion term. Inter-identities' interference was reduced to obtain a more expanded feature space. An unsupervised way to approach the challenge of visual ambiguities was presented in [187]. They have obtained the content and context from initial ranked lists and then removed the visual ambiguities from them to have the optimal feature representation and hence the final new rank list that was better than the initial rank list as tested on publically available datasets.

A view-invariant video-based few-shot learning method was proposed in [188]. Proposed method was developed on variational recurrent neural network that was trained adversarially to produce view-invariant feature space for matching persons. Another approach to mitigate the challenges *i.e.* viewpoint variation, low resolution and pose was presented in [189]. A lightweight and labelled part segmentation head was added to the backbone of re-id during training process to obtain diverse set of features hence resulted in improved performance of re-id.

In [190] an unsupervised cross-camera learning mechanism was adopted to achieve the generalized results. An unsupervised style transfer model was responsible to generate the images with transferred-style and different camera styles. Similarity was measured among generated and the original image. Similar images were then grouped together to form one cluster using iterative approach. A regularization term was also invoked to balance the cluster distribution. Results on benchmark datasets demonstrates the superior performance of the proposed approach. An asymmetric metric learning method was proposed in [191] to alleviate the view-bias problem. It was a two stream deep neural network that jointly learns the view and feature specific transformations. Clustering based deep asymmetric metric learning method was adopted to make the solution scalable.

A novel descriptor was presented in [192] for effective feature representation and metric learning. In order to have generalized multi-view the descriptor captures the structural information, analyzes and maximized the horizontal occurrences of multi-granularity to extract the rich feature representation even in case of drastic change in viewpoint. Besides the descriptor a metric learning method was also presented that jointly learns the multiple view-specific linear transforms to obtain the robust features. Descriptor and the metric learning method were then jointly evaluated on publicly available datasets to prove the effectiveness of the approach. A pragmatic semi-supervised framework was introduced in [193] to address the issue of view-specific biases. Framework learns the view-specific projections against each view. Only limited labeled data was used for training purpose. To boost the performance a re-ranking strategy was also introduced in the paper that measures the similarity among probe and gallery images and re-rank them based on their overlapping ratio. Framework has yielded superior performance when tested on mostly used re-id datasets in both supervised and semi-supervised way.

An adaptive multi-projection metric-learning method was introduced in [194] to handle the inconsistencies among different camera views. Proposed metric learning method jointly learns the different camera projections into a common feature space. Proposed method successfully adopts the newly added camera projections without updating the existing projection matrices. Notable improvements were observed when applied on major re-id datasets. A view-invariant subspace was learnt in [195] using adversarial approach. Specifically, coupled asymmetric mapping was learnt. View discrepancy was resolved by optimizing by cross-entropy view-specific objective. A similarity discriminator was introduced to determine the similarity value to distinguish the negative and positive pairs. To handle the imbalance of identity pairs caused due to most difficult samples adaptive weighing was also implemented.

Another deep multi-view feature learning method presented in [196]. Proposed metric learning based method exploits the fusion of handcrafted and deep learning features to produce the discriminative feature representation. An unsupervised framework for video-based person re-id was presented in [197]. Relation of frame with its first neighbour was explored to form the group in each camera. Cross-view matching strategy then find the matching relationship among them. And finally metric model for each camera pair was learnt in a progressive manner.

An approach was designed in [198] to address the variations exist in same video. Mainly a new loss term was defined comprised of intra-video loss and Siamese loss. Intra-video loss make the video more clustered by using the mean-body of each camera viewpoint. And Siamese loss placed the wrong matches more apart. Generalization capability of the model was increased as network was trained in iterative manner and hence the mean-body weights were updated accordingly. In [199] data discrepancy among multiple views was minimized using the proposed multi-level learning framework in iterative manner. Synthetic data was generated using already available grouping information and this data was then viewed as transitional state among original camera views. Gallery and probe images were moved into common subspace in a progressive manner to perform the matching step.

A memory module was proposed in [200] with the purpose to make the system invariant in terms of camera viewpoint and neighbourhood changes. They have used specifically unlabelled dataset to learn unsupervised discriminative representation and to make model domain invariance. ResNet-50 [127] as backbone with pre-trained weights on ImageNet [40]. An unsupervised approach to learn asymmetric learning of cross-view person images was presented in [201]. For each view model learns the specific projection, based on asymmetric clustering. In order to achieve the better matching performance, model finds the shared space with low view-specific bias.

Zimo Liu *et al.* [108] proposed an unsupervised tracklet based framework to learn cross-view discriminative features. They have used tracklet as query in search for nearest neighbour best match after possible iterations until the best match found. In order to reduce the impact of false positive matches they have employed hard negative mining. KNN searching was then repeated in a reverse manner to ensure the best match found. Best query matches at initial stage and at reverse stage are then used collectively to update the model.

In [202] an unsupervised re-id framework to cope the challenge of viewpoint variations. Deep Clustering-based Asymmetric Metric Learning (DECAMEL) learns an initial asymmetric metric using a linear unsupervised model *i.e.* CAMEL. It then embeds the learned metric into deep network by jointly learning metric and features in an end to end manner. Afore mentioned steps were based on asymmetric metric clustering. Novel loss function was then applied to achieve the best results. Sub-optimality in the results was obtained on the basis of separation of metric and feature learning. Due to learning of better cross-view clusters in the shared space better cross-view matching performance was achieved.

5.6.2. Attention-based Approaches:

Meng Zhenget *al.* [203] has proposed a Siamese architecture to address the challenge of viewpoint. In proposed joint learning end-to-end architecture a flexible attention mechanism was introduced to achieve the attention consistency among the images of same person. Lei Zhang *et al.* in [204] to address cross-view challenges. In the proposed end-to-end framework view-invariant features were learnt and comprised of three components *i.e.* adversarial learning, drawing same features towards center and SIFT guidance. To improve the integration of these components attention mechanism was also adopted to produce the superior results.

In [205] attention aligned network was presented that focuses on foreground information using channel wise multi-scale attention aware mechanism that had helped in learning the invariant views obtained from different cameras. To increase the capability of feature learning an improved triplet loss was also presented. Resulted in improved results by maximizing the inter-class distance and minimizing the intra-class distance.

5.6.3. GAN-based Approaches:

To learn the view-invariant features GAN and another contrastive learning module was combined into one training framework in [206]. Novel views are generated using mesh based view generator. Proposed method was flexible as the model does not rely on labeled source domain. Improved results are obtained as compared to existing fully unsupervised and unsupervised approaches.

Table 10: Results obtained on viewpoint variation challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	CAMEL, 2017, [201]	Market-1501	54.5/26.3	No
		CUHK-01	61.9/57.3	
		CUHK-03	39.4/31.9	
		VIPeR	30.9/-	

2	LOMO+GOG, 2017, [108]	MARS iLIDS PRID-2011	23.9/– 41.7/– 80.9/–	No
3	DVDL, 2015, [174]	iLIDS PRID-2011	25.9/– 40.6/–	No
4	VKD, 2020, [176]	MARS DukeMTMC-VideoReId	88.74/82.22 95.01/93.41	Yes
5	GCL, 2021, [158]	Market-1501 DukeMTMC-ReId MSMT-17	90.5/75.4 81.9/67.6 54.4/29.7	Yes
6	PersonX, 2019, [173]	Market-1501 DukeMTMC-ReId	93.0/80.8 83.6/72.6	No
7	ECN, 2019, [200]	Market-1501 DukeMTMC-ReId MSMT-17	63.3/40.4 63.3/40.4 30.2/10.2	No
8	CASN, 2019, [203]	Market-1501 DukeMTMC-ReId CUHK-03(Labeled) CUHK-03(Detected)	94.4/82.8 87.7/73.7 73.7/68 71.5/64.4	No
9	LOMO-XQDA, 2015, [175]	CUHK-03(Labeled) CUHK-03(Detected) VIPeR	52.25/– 46.25/– 40/–	Yes
10	FC2, 2018, [184]	Market-1501 CUHK-03	48.06/– 37.41/–	No
11	ICV-ECCL, 2018, [181]	Market-1501 CUHK-01 CUHK-03 VIPeR	90.6/77.3 83.5/– 88.6/– 51.9/–	Yes
12	BINet, 2021, [185]	Market-1501 DukeMTMC-VideoReId CUHK-03(Labeled) CUHK-03(Detected) MSMT-17	95.3/88.7 91.4/81.3 73.6/72.5 72.3/69.8 76.1/52.8	No
13	SMC-ECD, 2018, [186]	Market-1501	80.31/59.68	No
14	DCIA, 2021, [170]	PRID-2011 VIPeR	32.5/– 64.78/–	No
15	VRNNs, 2020, [188]	MARS iLIDS	61.2/52.1 64.6/–	No
16	MGN, 2020, [189]	Market-1501 DukeMTMC-ReId CUHK-03 (Labeled) MSMT-17	95.8/88.7 90/79.9 78.8/74.4 84.0/62.4	No
17	UDA, 2020, [190]	Market-1501 DukeMTMC-ReId	73.3/38.0 56.1/30.6	No
18	VIH-ReID, 2015, [179]	VIPeR	21.4/–	No
19	CRAFT-MFA, 2018, [180]	Market-1501	77.0/50.3	No

			CUHK-01	74.5/–	
			CUHK-03	84.3/72.41	
			VIPeR	50.3/–	
20	DECAMEL, [202]	2020,	Market-1501	60.24/32.44	Yes
			CUHK-01	65.81/–	
			CUHK-03	38.27/–	
			MSMT-17	30.34/11.13	
21	CSPL, 2018, [177]		CUHK-01	72.02/–	No
			CUHK-03 (Labeled)	70.2/–	
			CUHK-03 (Detected)	66.8/–	
			VIPeR	51.3/–	
			PRID-2011	69.20/–	
22	DAM, 2019, [191]		MARS	74.65/–	No
			iLIDS	77.3/–	
			PRID-2011	87.0/–	
23	GMDA-RC, [192]	2018,	CUHK-01	75.69/–	No
			VIPeR	67.09/–	
24	VS-SSL, 2020, [193]		Market-1501	74.8/51.2	No
			CUHK-01	73.0/–	
			VIPeR	44.8/–	
25	CVDCA, 2016, [178]		CUHK-01	47.8/–	No
			VIPeR	47.78/–	
			PRID-2011	57.6/–	
26	VCFL, 2021, [204]		Market-1501	91.85/76.97	No
			DukeMTMC-ReID	82.68/65.68	
			CUHK-03	61.29/54.26	
27	MPML, 2019, [194]		Market-1501	55.61/27.58	No
			CUHK-01	64.98/–	
			VIPeR	44.72/–	
28	AANet, 2021, [205]		Market-1501	96.1/87.5	No
			DukeMTMC-ReID	90.2/79.5	
			CUHK-03 (Labeled)	82.7/76.7	
			CUHK-03 (Detected)	77.2/70.5	
29	AVA-ReID, [195]	2020,	Market-1501	88.6/73.1	No
			CUHK-03	86.7/83.8	
30	CCM, 2021, [197]		DukeMTMC-VideoReID	66.0/42.3	No
31	RCN-ReID, [198]	2019,	MARS	48.0/–	No
			iLIDS	65.0/–	
32	MLCPL, 2018, [199]		CUHK-01	34.05/–	No

5.7. Low Resolution

Mostly surveillance cameras are not capable to capture the high resolution images because of low resolution of cameras and vast distance between person and camera as shown Fig. 1. Low resolution probe images and high resolution gallery images makes the re-id more challenging. Due to variations in pose, illumination and resolution appearance of same person may look very different in non-overlapping cameras. Technique of image restoration unable to produce efficient results in real time scenarios leveraging the domain gap in self-supervised manner *i.e.* without using extra cost of labelling. Therefore, it is important to address this problem and existing approaches address it by using appearance based methods.

These methods extract feature representation that contain high inter-class disparity between different subjects and low intra-class disparity for same subject. The intra-class disparity is often extensive than the inter-class disparity due to the significant appearance change across different cameras. Therefore, accurate matching become difficult.

Additionally, illumination variance is one of the important challenge to address. As the light conditions play the most important role to match the query person from a large set of images captured by the non-overlapping cameras. Sometimes the color of body parts of a particular person perceives so different due to complex illumination variations as shown in Fig. 1.

Changes in visual appearance caused due to light variation is another challenge in person re-id. As variant light conditions also causes the change in pixel values. Lighting conditions significantly affects the performance of low level features such as texture and color. Color information of pedestrian like clothes become unidentifiable when lighting conditions changes (with or without lighting). Existing approaches are based on learning the depth information, color based features and joint representation learning. To extract the depth information using depth cameras such as Kinect is not a difficult task. Kinect obtain the depth value (distance to the camera) of each pixel by infrared, regardless of subject color and illumination. But there are several limitations, depth images captured by depth cameras changes when the view point of person changes across camera views. To develop a model which efficiently learn the discriminative low level features is still challenging task in person re-identification. Fig. ?? shows the state of the art results on respective datasets achieved so far on illumination challenge.

Fig. 13 shows the progress of available papers on low resolution and illumination variance between 2015 to 2021.

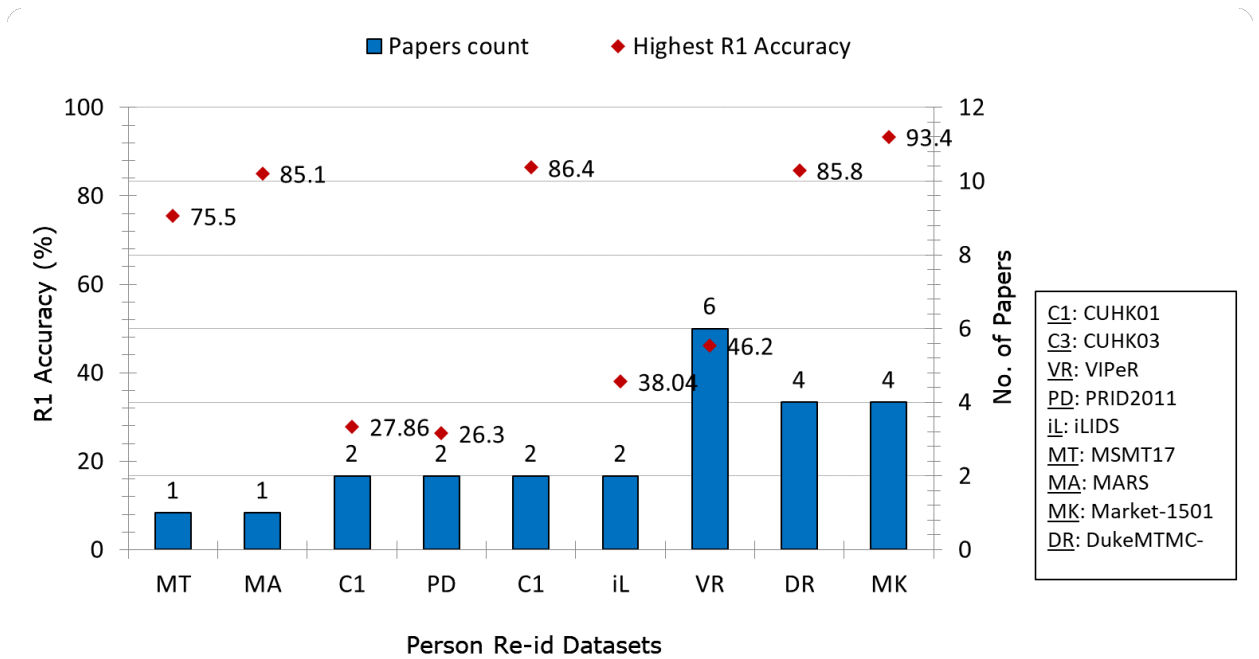


Figure 13: Progress on the challenges of low resolution and illumination variance for person re-id benchmarks

5.7.1. CNN-based Approaches:

A supervised framework to resolve the challenge of resolution was proposed in [207]. They have combined embedding of multiple layers into single layer. In order to achieve efficient results lower layers with higher resolution are combined with higher layers having semantic information. Xiang li *et al.* [208] presented an attempt to resolve the low resolution that was beyond the scope of re-scaling and interpolation mechanism. They have presented the image features that are at same scale in a latent space and then performed distance metric modeling at each scale, this resulted in formation of shared space among low resolution and normal images of same person to obtain the effective results.

In [209] issue of no availability of super resolution images was challenged. In the presented procedure pair of high resolution and low resolution dictionaries and mappings functions are learned during training, due to this learned

dictionary and mapping function low resolution images would be converted to high resolution discriminant features. Ke Han *et al.* in [210] challenged the issue of low resolution. Presented model predicts and recovers the content aware details. Using self supervised approach the model also assigns soft labels that are automatically generated. Important thing is unique labels are assigned according to the optimal scale factor, it is to recover the issue of lacking ground truths. The probability of generated labels indicates the optimality of assigned scale level. Hence increasing the level of confidence against optimal scale for optimal resolution with context aware prediction.

In [211] an end-to-end adaptive feature fusion framework was proposed that has proved effective in handling resolution at different recovered body regions and at multiple scales. An adaptive feature integration module balances the relative importance of super resolved image content. In this manner adaptive weights were assigned to input features with super resolution. An extended version of [209] was presented in [212] to improve the results on resolution and they have showed analysis results on two more datasets with improved approach. Now the structure was supporting multi-view. They have used to learn different mappings to convert low resolution images into discriminative high resolution features. In [213] resolution-aware framework was proposed. They have used the knowledge transfer technique to minimize the variations in resolution of images. Teacher knowledge was exploited and transferred to the low resolution student network to narrow down the resolution differences.

Mainly, challenge of illumination, camera-viewpoint and resolution were focused for adaptive cross-domain person re-id. In [214] illumination/lighting change was challenged. Reinforced temporal attention based end-to-end framework was proposed that was implemented at each frame level features to extract the temporal information accurately for depth based person re-id. LSTM was used after obtained features from CNN. LSTM models/learns the short term temporal dynamics/changes. Proposed approach is also useful for person with unseen clothes.

For metric learning a multi-modality approach was presented in [215], it learns the changes in illumination via shift-invariant property. Proposed model also learn the sub-metric against each modality to reduce the role of bias in a global sense. Approach was validated on multiple datasets.

In [216] illumination impact was studied and proposed a framework to resolve. It transform the pixels to invariant color space. Proposed framework learns the patterns and structures inherent in image pixels. It jointly learns the encoding and transformations among pixels pairs belonging to image pairs. This auto-encoder based approach transform the 3D-RGB pixel values to higher dimensional space and encode them using dictionary. Then mapping of pixels to invariant space is performed. These encoded pixel values are pooled over a particular regions and then integrated (or concatenated) to formulate the final representation. For higher level complementary feature learning, this framework can also be extended to multi-layers.

5.7.2. GAN-based Approaches:

Cross resolution problem was considered in [217]. In the proposed methodology existing GAN architecture was improved using high resolution images in an end-to-end fashion. Problem of low resolution was resolved by forming the association between super resolution images and re-id task. The formulation was due to parameterized sharing while training in end-to-end manner.

The feature mismatch problem caused due to the variations in resolution and illumination might lead to poor representation learning and was resolved in [218]. In this self-supervised strategy identity related information was extracted to resolve the challenge of degradation occurs in real time. No extra computation was used instead model works on low resolution images. The approach was based on GAN, in which images were generated by switching the content of real world images with generated images. In this way domain gap between gallery images and real world images was reduced. Model effectively preserves the features related to identity information and remove the features that were related to degradation to achieve the effective results.

Table 11: Results obtained on low resolution and illumination variance challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	JUDEA, 2015, [208]	VIPeR	25.87/–	No
2	PRI, 2020, [210]	Market-1501	86.9/–	No
		DukeMTMC-ReId	82.1/–	
3	INTACT, 2020, [217]	Market-1501	88.1/–	No

		DukeMTMC-ReId	81.2/–	
		CUHK-03	86.4/–	
		VIPeR	46.2/–	
4	DaRe, 2018, [207]	Market-1501	90.9/86.7	Yes
		DukeMTMC-ReId	84.4/80.0	
		CUHK-03 (Labeled)	73.8/74.7	
		CUHK-03 (Detected)	70.6/71.6	
		MARS	85.1/81.9	
5	SLD2L, 2015, [209]	CUHK-01	24.48/–	No
		VIPeR	16.86/–	
		iLIDS	33.33/–	
		PRID-2011	22.6/–	
6	MCSLD2L, 2017, [212]	CUHK-01	27.86/–	No
		VIPeR	20.79/–	
		iLIDS	38.04/–	
		PRID-2011	26.30/–	
7	RKD, 2021, [213]	Market-1501	93.4/83.7	No
		DukeMTMC-ReId	85.8/73.0	
8	DI-REID, 2020, [218]	MSMT-17	75.5/–	No
9	JLCF, 2016, [216]	VIPeR	26.27/–	No
10	M3L, 2017, [215]	VIPeR	30.22/–	No

5.8. Cross-Domain/ Generalization

Person re-id models produce improved results when trained and tested on same dataset but perform poor when tested on different dataset due to different scenarios *e.g.* changes in viewpoint, place, background, resolution and different visual appearance. For this, recently clustering, domain adaptation and image to image translation based approaches reported state of the art results. One such approach is shown in Fig. 14. Aims of image to image translation is to develop a mapping function between two domains and it required paired training data which is difficult to manage. For domain adaptation, labelling a dataset is expensive and time consuming task. Due to these factors improving generalization is still challenging in person re-identification. Progress of Papers that has addressed the generalization challenge in top conferences and journals is shown in Fig. 15

5.8.1. CNN-based Approaches:

Changes in background and poses results in challenge of intra-class variation and that is addressed in [219]. In the proposed discriminative module primary features are learned along with generative module to generalize the data well. In [128] semantic segmentation with simple base line was used to address the challenge of local feature extraction from various human body parts that is an alternative to bounding box approach. In this framework, in order to integrate the body part classification, the identity attribute features was presented. ResNet-50 [127] along with pre-trained ImageNet weights was used to resolve the partly attention-based challenge. Ruijie Quan *et al.* [220] proposed an architecture that automatically searches for appropriate CNN that is suitable for efficient person re-id eliminating human effort. Body part information was used to capture the structure information of human body parts. Structure information was then embedded into flexible part-aware module to achieve state of the art results.

Poor generalization of the model may lead to domain gaps that are needed to be covered. One of such approach to tackle this challenge was given in [221] using unsupervised approach. Specifically a style normalization and restitution module was introduced that filters out the variations in style and color using instance normalization. Features related to identity are then refined and added to the network to achieve valuable discrimination results on unseen data as well. For better feature disentanglement dual loss was also used in the proposed strategy. In [153] Houjing Huang *et al.* designed a cross-domain adaptive model without identity labels of target-domain. Part Aligned Pooling (PAP) was introduced to enhance alignment and hence results in model generalization. In the approach it is verified that part-alignment plays a vital role in achieving cross-domain generalization of re-id model. They have argued that training unlabelled target domain with part segmentation and training re-id on source domain, is an effective way to achieve cross-domain generalization or domain adaption.

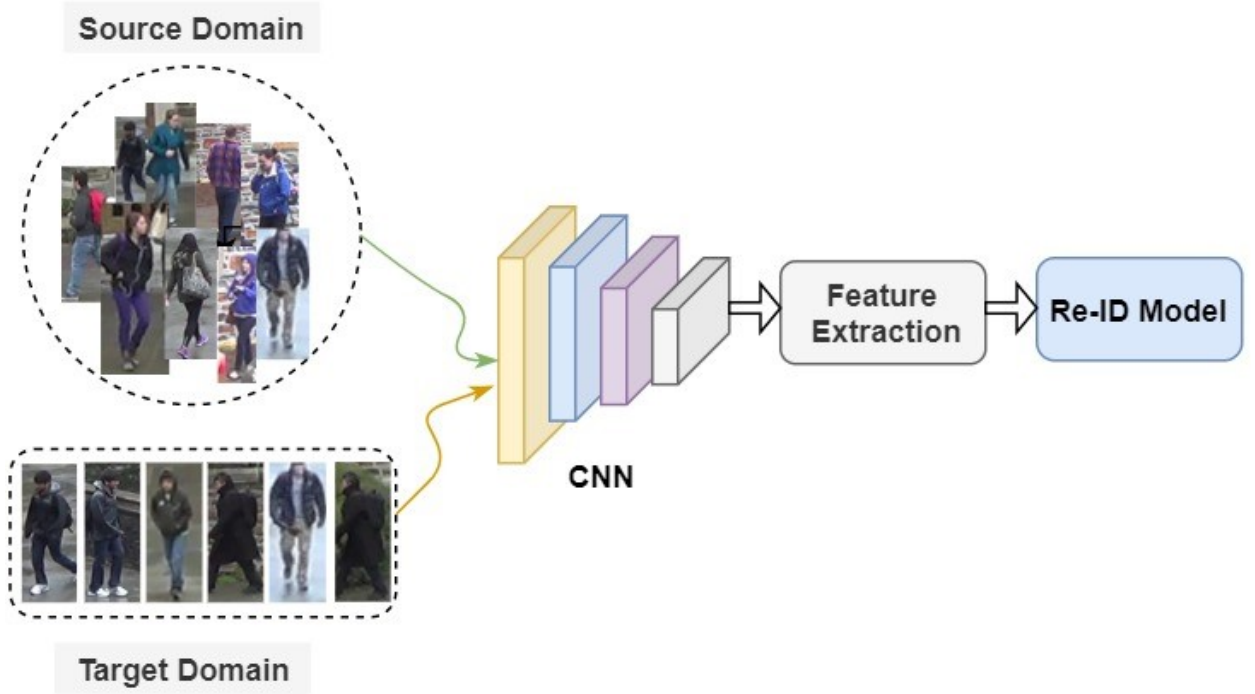


Figure 14: An approach to achieve generalization.

Weihua Chen *et al.* in [222] presented novel quadruplet ranking loss that was based on triplet loss; it was to increase the generality of the procedure for person re-id. It effectively achieves the capability to handle the intra and inter class variations. A quadruplet deep network was also presented that makes 4 hard samples of dataset and then proposed loss was applied to obtain the effective cross-domain results. To make the model domain invariant an efficient procedure was presented in [223]. At first stage, for strong baseline they have merged multiple person re-id datasets and trained a CNN using single softmax loss. Standard dropout layer was replaced with domain guided dropout layer for training few more epochs, resulted in effective learning of weights that are more effective for each domain.

In [224] domain adaptive person re-id was proposed using unsupervised target domain. Multiple expert brainstorming networks were learned with multiple architectures for optimal re-id. In proposed network models were first pre-trained on multiple expert models. To adjust the heterogeneity of multiple models a regularization scheme was added that modulates the expert models according to their feature distribution in the target domain. In this way significant performance was gained due to the increase in discrimination capability of re-id model. Another framework to achieve generalized results in unseen scenarios as well was proposed in [225]. In order to reduce the domain gap in the presented mechanism all camera images were assigned the common subspace. This resulted in achieving the generalized results even on unseen images.

In another framework presented in [226] interpolation mechanism was used to obtain cross-domain generality. Vanilla neighbourhood approach was improved by restricting the camera-aware manner. In [227] a joint learning framework was presented for generalized person re-id. Representation space was purified in a way that it can only learn id related feature space. Disentangling module encodes cross-domain images into appearance space and structure space that was shared for reduced domain gap. Recently a new dataset Person30K is proposed in [228] to support generalization. Proposed DMG-Net strengthens the model to attain generalization on unseen data without the need to fine tune further. Specifically, feature based centers were computed that represent identities and then matching was done between support centers and query samples.

In [229] Meta batch instance normalization (MetaBIN) technique was used in an unsupervised manner to achieve the generalization on unseen person re-identification domains. The MetaBIN approach prevents the model from over-fitting. It does so by learning the scenarios before hand in a phase called meta-learning. Moreover meta-train loss accompanied by a cyclic inner-updating manner helped in boosting the generalization capability. It is noted

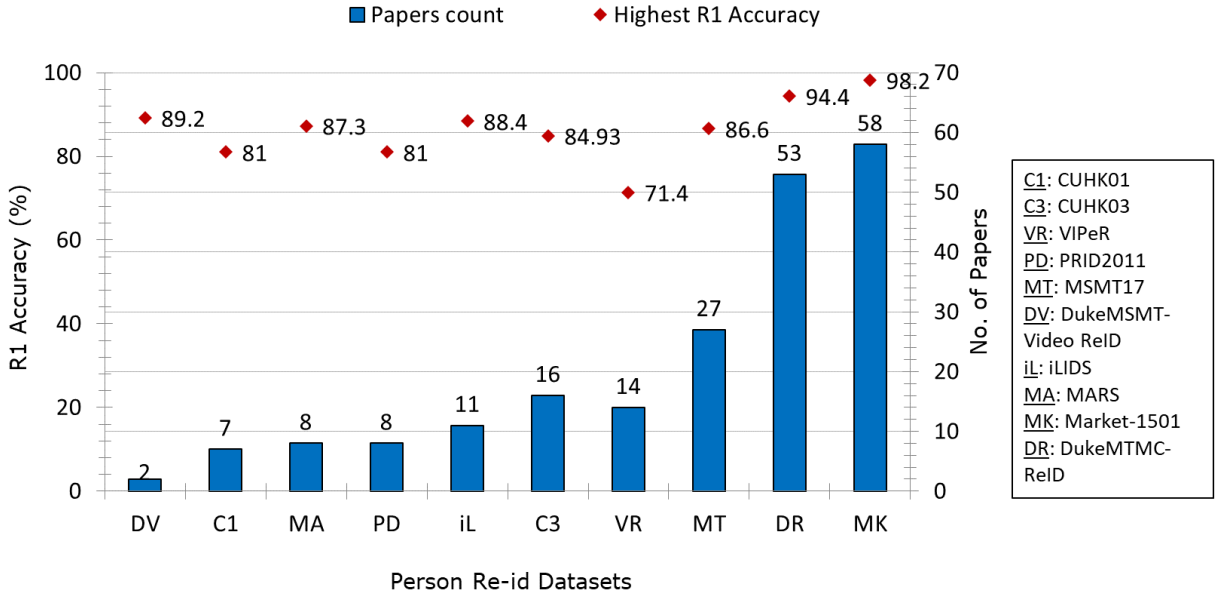


Figure 15: Progress on the challenge of cross-domain generalization for person re-id benchmarks

that no additional augmentation technique used that makes the model more complicated. Another method was proposed in [230] named as relevance-aware mixture of experts (RaMoE). In the method voting based technique was used to dynamically integrate the diverse characteristics of source domains that eventually helped in increasing the generalization capability of the model.

Anguo Zhang *et al.* in [231] an issue of intra-class variation was challenged. Second order information bottleneck was integrated into the network. It optimizes the network at inference time resulted in reduction in computation overhead. Proposed framework has reduced the impact of intra-class variations and superior performance was achieved. To generalize the results on unseen domains experiments on synthetic data were presented in [232]. Usage of unreal data have significantly reduced the cost of both training and testing and have also improved the results via using the transfer knowledge technique on real datasets. The approach has also provided the basis for unsupervised domain adaptations. The pre-trained model obtained via proposed technique can be easily plugged into existing methods to achieve better results.

An optimization strategy was presented in [233] to achieve the generalization. Specifically label banks were formed in hierarchical manner, mini-batches were updated using iterative approach. Hierarchical clusters form and split on the basis of computed or propagated label information. Effective labeling resulted in achieving the generalization of the proposed model. [234] generalization was improved that was based on their previous work [235]. An additional term was added in the quartet loss that ensures the distance between matching pairs should be less than specified threshold, that resulted in less intra-class distance and more than inter-class distance with different probe images. It also ensured the closeness of intra-class features to each other. Moreover, model has also optimised both identification and classification tasks.

In order to reduce the difference in domains due to background inferences and to attain generalization, an unsupervised approach was presented in [120]. They have suppressed the background information to have the focus on learning the person's information present in the foreground. They have used human body information and associated ID related information present in the environment. Model keeps on updating using the virtual labels of target domain. In an effort to learn the efficient feature embedding a joint learning framework was presented in [236]. In order to predict a non-ambiguous label information subgroups were formed by merging neighbours. For higher subgroup influence similarity-aggregating loss was introduced that has pulled the similar embedding closer to each other.

Using multi-tasking reinforcement learning both spatial and temporal embedding were learnt by a framework presented in [237]. The framework was comprised of two branches. First branch determines the optimal spatial region

along with temporal range of frames to obtain the contextual information. While second branch has focused on learning the identity information based on the information obtained from the first branch. Collaborative interaction among pedestrians and the context were exploited to achieve outstanding performance. Another approach to mitigate the challenges of unsupervised domain adaptation was proposed in [238]. In the proposed dual-refinement approach pseudo labels were refined at off-line clustering while feature refinement was performed at on-line training phase. This dual-refinement has minimized the influence of noisy labels and hence superior results were obtained. A model to learn the domain invariant features was designed in [239]. Adversarial auto-encoder was the model base. And maximum mean discrepancy measure was used to align the learned distribution across variant domains.

A 3D guided network was developed in [240]. In order to support the formulation on unknown domains as well the proposed model have the ability to transfer knowledge from domain training data. They have performed image-to-image translation to form the possible synthetic poses and viewpoints of a person while preserving the identity information. In a way higher accuracy results were achieved even on unknown domains. Aiming to enhance the generalization with low computational complexity, another framework was presented in [241] comprised of two related modules *i.e.* patch-based metric learning and local salience learning. CNN was used to extract the features of person patches then just two positive patch-pairs were selected to learn the patch based metric matrices. k-means clustering was then used as a salience learning algorithm to learn the patch weights. Finally, on the basis of these learned weights patch-wise similarity score was computed.

An approach to address cross-dataset person re-id was developed in [242]. In the proposed dictionary learning approach visual-attribute embedding was learned and then transferred from source to target domain. Finally, pseudo-labels were fine-tuned by choosing few samples from the target domain to produce outperforming results. The challenge of part-based learning was addressed in a model named PatchNet [243]. In order to learn the patch features discriminative model measures similarity between images in unsupervised manner. In [244] a deep representation learning procedure named as Part Loss Network (PL-Net) was proposed. Part loss network minimizes the empirical classification risk (ECR) on training images and representation learning risk (RLR) on unseen images. RLR was evaluated via part loss. It was evaluated for each body part separately. This part loss computation helped in gaining the discriminative power on unseen images as compared to traditional global classification loss.

A semi-supervised approach based on deep attribute learning framework was presented in [245] to achieve cross-domain generalization. At first stage fully supervised training is performed, resulted in learned attribute labels against target dataset. Second stage involved the fine tuning of learned labels on the basis of attribute triplet loss. At every iteration, attribute triplet loss ensures that similar attributes belong to same person. Second stage predicts attributes for target dataset, these are then combined with initial labels obtained at first stage and then fine tuned again. At final stage, more discriminative deep attributes were obtained and that have performed impressively well. The approach was semi-supervised because it includes one dataset with labels and other without any attribute label.

In [246] context representation was learned and transferred to target re-id dataset via Bayesian adaptation. Proposed approach can handle both labelled and weakly labelled data in the training process. Learned representation was view invariant but person variant, that is an ideal case for efficient person re-id. To make re-id task more scalable a domain transfer unsupervised technique was presented in [247]. In the approach they have used target data with no labelled matching pairs. Model transfers the view invariant representation from labelled datasets to an unlabelled target dataset. Dictionary learning was adopted assuming that multiple appearances of a person can be represented as linear combination of vectors.

An unsupervised approach to achieve the effective outcome in real world scenarios was presented in [248]. Domain gaps were managed using Gaussian distribution. Positive and negative samples were pushed apart using unsupervised setting. A momentum update mechanism was applied to clearly separate the negative and positive samples distribution. Hence significant improvement over baselines was achieved. In [249] potential of cross domain model was explored for real world application of person re-id systems. System was developed using unsupervised domain adaptation. Proposed model has adaptively learned the credible samples for training by avoiding noisy labels. An instance margin loss was introduced that has increased the margin of instance. Resulted in increased instance level discrimination of learned features.

Djebril Mekhazni *et al.* in [250] presented an unsupervised scheme for unseen domain adaptation for person re-id. Pair-wise distances were optimized by using gradient descent and comparatively smaller batches. They have introduced novel dissimilarity based discrepancy loss that made the source and target distribution similar to make the results more efficient. Multi-source domain training to make the model generalize on multiple unseen domain was a way adopted in [251]. For this purpose, there comes a problem of domain gaps. To resolve such issues two strategies were adopted. First

was domain-specific batch normalization and second was the infusion of multi-domain information. These strategies had helped in reducing the domain distances via fusing the features of multiple domains. Proposed techniques have provided the comparable results to supervised approaches. Attained results were obtained without using any post-processing techniques. In order to address the limited scale issue of existing re-id datasets, a new dataset named as Large scale Unsupervised Person re-id (LUPerson) was introduced in [252]. They have investigated the impact of factors in learning *i.e.* data augmentation and temperature usage in a contrastive learning framework.

Degraded label assignment caused due to distribution discrepancy among cameras was a key issue and addressed in [253]. To obtain the possible accuracy on pseudo-label's generation the proposed framework was divided into two modules *i.e.* inter-camera and intra-camera computation. In inter-camera computation a new feature vector was computed for different cameras. This new feature vector results in obtaining the reliable pseudo-labels and hence removes the issue of distribution discrepancy among different cameras. While in intra-camera computation similarity measure was computed by comparing features obtained from CNN. Then the model was trained in two stages *i.e.* inter-camera and intra-camera pseudo-labels to obtain the outperforming results. In [254] problem of discriminative learning with unlabeled data for unsupervised person re-id is targeted. A Dynamic and Symmetric Cross-Entropy (DSCE) loss was introduced to cope the challenge of negative effect caused due to noisy labels. And for handling the camera shift issue a meta-learning (MetaCam) technique that splits the training data into meta-train and meta-test was used. MetaCam resulted in achieving the interactive gradient impact for each meta set that would enforce the model to learn the camera-invariant features for better results.

To tackle the challenge of pseudo label noise an approach was presented in [255]. They have refined the pseudo labels based on clustering consensus. Pseudo labels were refined dynamically with temporally generated ensembles pseudo labels. This refinery approach can also be integrated into existing clustering based methods to obtain the simple yet effective results. A method with only few labeled information was presented in [256]. They have presented a multi-domain generalization framework that can perform well on unseen domains with no need to train a new model. Specifically, a meta learning strategy was adopted to perform the like train-test process resulted in learning more generalized model. A memory based generalization loss and a meta batch normalization layer was also added to diversify the advantages of meta-learning. In [257] a Group-aware label transfer algorithm was proposed. It promotes the pseudo-labels via online interaction. It not only uses the pseudo-labels but at the same time it also refines them based on an online clustering algorithm. The effectiveness of the approach was tested on large-scale re-id datasets and hence it has helped to reduce the gap among supervised and unsupervised approaches.

In order to bridge the domain gap collection of unsupervised schemes were adopted in [258] that have resulted in learning the unique representation of features. First a clustering algorithm optimizes the features in iterative manner to ensure that learnt features are noise free. Second was the progressive domain adaptation and third was the Fourier augmentation in which extra constraints were deployed to increase the chances of class separability. In [259] patch-wise features were learnt from unlabelled patches of person's image. A novel loss function was then used that guides the model to mine the discriminative information. Same features were pulled closer while features belonging to different instances were pulled away depending on the computation of loss function.

Yan Bai *et al.* in [260] a reliable approach to generate pseudo labels was presented for unsupervised domain adaptation. In the proposed hierarchical scheme graph convolutions were used to learn the complicated structure of each cluster. Connection among samples were estimated in a hierarchical way and then refined progressively. In video person re-id [261] a joint global video and local frame information was considered to obtain the diverse information for better estimation of pseudo labels. Moreover novel loss term induce the model to not focus on undesirable factors of identity. A dynamic strategy was also adopted to choose the pseudo label with higher confidence score that keeps on updating during training process until the higher confidence score met.

An end-to-end self-supervised learning algorithm was proposed in [262] for unsupervised person re-id. Domain discrepancy was minimized using agent learning mechanism. Due to learnt discriminative representation model has performed well on unseen domains as well as compared to existing unsupervised approaches. Hang Zhang *et al.* has presented an unsupervised view-invariant approach to tackle the person re-id at multi-scale levels in [172]. Proposed framework was self-trained and can be implemented on unseen domains as well. Local and global representations of pedestrian images were learnt. And then used as pseudo-labels that were improved progressively using iterative approach. For unsupervised domain adaptation a multi-loss optimization learning model was proposed in [263]. They have estimated the pseudo labels via clustering mechanism applied in a supervised way. For similarity and adversarial learning two losses were introduced that had helped in model optimization. Together these loss terms has also benefited in exploring the intra-domain relation to evaluate and estimate the improved final outcome.

In [264] discriminative representation was obtained from tracklet data in an end-to-end formulation. Proposed framework learned the discrimination and association within-camera and cross-camera respectively. Superior performance on eight benchmark datasets proves the effectiveness of the proposed framework. In [265] a self-supervised learning algorithm was proposed to achieve generalization. In order to bridge the gap among source and target attribute-identity embedding were used as a base to optimize the model. A prediction-training cycle was also implemented with a purpose to fine-tune the model variables so that the model become more adaptive to target domain. In [17] video tracklets were used to achieve the generalized results at large scale, hence proposed new ideal dataset named MARS. They employed motion features and CNNs to learn discriminative embedding. According to their results motion features are less effective in real scenarios because of complex challenges like occlusion, pose and background clutter. However CNN based features remarkably performed well as there we have large training data hence a good generalization ability was achieved.

Challenge of cross-view and intra-bag alignment was addressed using weakly supervised approach [266]. To reduce the labelling cost presence of humans was labelled only however the information of what and where was not considered. Density based clustering is a technique adopted by [267] to form a model that is domain adaptive. Proposed model was made discriminative in the target domain using GAN's min-max strategy. Sample clusters were first predicted from the target domain, and then sample features were extracted using re-id model. The model was already pre-trained on the source domain. Cluster formation was then improved via iterative process. Feature encoder increases the inter-class distance and decreases the intra-class distance. Image generator and feature encoder competes in adversarial min-max manner and hence helped to optimize the models effectively.

Liao *et al.* in [268] also considered the generalization as a challenge and proposed a framework to resolve it. In order to bridge the domain-gap 'divide-and-conquer' strategy was adapted for factor-wise style transfer [73] based on CycleGAN [269]. An unsupervised domain-to-domain translation method was presented in [270]. The method keeps the pedestrian identity information and have used maximum mean discrepancy as a base to pull the alike distribution closer. They have used CycleGAN to transfer the label information to unlabeled domain. Weijian Deng *et al.* [271] has improved the baseline of "learning via translation", in their proposed work they have tried to generalize the domain by preserving the similarity of image and dissimilarity of domain in unsupervised manner. Their proposed framework was based on Siamese Structure and Cycle-GAN.

In [272] generalization was targeted using lifelong learning scenario. In the proposed framework learn-able knowledge graph was maintained that updates the previously learned knowledge in an adaptive manner. In order to improve the generalization learned knowledge was transferred on unseen domains. Promising improvement were reported over existing SOTA results. A self-training scheme for optimized unsupervised domain adaptation was designed in [273]. An encoder was trained on the basis of guessed labels for unlabeled target data to achieve the effective results. New open environment larger dataset was introduced in [10] named "Market" to achieve generalized results due to its large scale. They have also presented an unsupervised Bag of Words descriptor that takes re-id task as image search and this is beyond the scope of this paper and hence not further discussed.

5.8.2. Attention-based Approaches:

A self-critical attention based learning mechanism was proposed in [48]. In the proposed design unified modules self-critic and self-correctness are introduced to guide the attention agent learn the correct attention maps based on information provided by critic module. In [274] high-order attention module was proposed to resolve the challenge of part-based modelling of pedestrians in person re-id task to generate more discriminative and powerful attention proposals. High order relationship among human parts was obtained using polynomial predictor of high-order. In this way discriminative attention maps are obtained with subtle differences. Proposed model works well on unseen person images as well due to learning at multiple diverse levels so that attention of all sides could be preserved.

In [275] classification based attribute aware re-id approach was proposed. In the proposed model attention mechanism was used to identify the specified body parts in a unified learning framework. Architecture integrates the identity information with attribute features and body parts. In this manner discriminative feature space was learned and model becomes more generic. An end-to-end supervised approach [276] to make the model context free was based on attention mechanism. In the proposed Siamese network intra sequence and inter-sequence attention mechanism was used for feature refinement and alignment accordingly. A novel cross-correlated attention module was presented in [277] that effectively learns the discriminative representation that was based on inherent spatial relation of different regions of a person image.

In [278] large-scale re-id scenarios are targeted. A novel harmonious attention network was deployed that jointly learns the attention based pixel representation of soft and hard regions. Hence resulted in having more discriminative features that has helped in more efficient re-id searching and matching. An attentional aggregation formulation was designed in [279] to handle the changing representation of an identity in a query image. Proposed scheme flexibly incorporates the similarity metrics along with multiple representations. Videos captured from different cameras might end up a video with different camera views, and this view variation is challenge handled in [280]. Their end-to-end framework accounts the inter-dependencies among video sequences. They used Recurrent Convolutional Network to extract features and then learns similarity among them. These similarity scores are then used to form the attention network in spatial as well as temporal dimensions. Obtained attention vectors are forwarded for pooling. At final stage, Siamese structure is deployed at attention vectors to make the solution more generalised.

Table 12: Results obtained on Cross-domain/Generalization challenge against each dataset. Results in bold are the highest.

Sr.No	Paper	Dataset	R1/mAP	Code availability
1	CCL-PDA-FA, 2021, [258]	Market-1501	94.2/83.4	No
		DukeMTMC-ReId	83.5/70.8	
		MSMT-17	66.6/36.3	
2	OPLG, 2021, [233]	Market-1501	91.5/80.0	No
		DukeMTMC-ReId	82.2/70.1	
		MSMT-17	56.1/29.3	
3	SCAL, 2019, [48]	Market-1501	95.8/89.3	No
		DukeMTMC-ReId	88.9/79.1	
		CUHK-03 (Labeled)	74.8/72.3	
		CUHK-03 (Detected)	71.1/68.6	
4	NAS, 2019, [220]	Market-1501	95.4/94.2	No
		CUHK-03 (labeled)	77.9/73	
		CUHK-03 (Detected)	73.3/69.3	
		MSMT-17	78.2/52.5	
5	MHN, 2019, [274]	Market-1501	95.1/85.0	Yes
		DukeMTMC-ReID	89.1/77.2	
		CUHK-03	71.7/65.4	
6	ASTPN, 2017, [280]	MARS	44.0/–	Yes
		PRID-2011	30.0/–	
7	DPM, 2015, [?]]	Market-1501	42.64/19.47	No
		CUHK-03	22.95/22.7	
		VIPeR	21.74/26.55	
8	MEB-Net, 2020, [224]	Market-1501	89.9/76.0	Yes
		DukeMTMC-ReID	79.6/66.1	
9	GDS, 2020, [248]	Market-1501	81.1/61.2	No
		DukeMTMC-ReID	73.1/55.1	
10	DCML, 2020, [249]	Market-1501	88.2/72.3	No
		DukeMTMC-ReId	79.3/63.5	
11	NRMT, 2020, [120]	Market-1501	87.8/71.7	No
		DukeMTMC-ReId	77.8/62.2	
		MSMT-17	45.2/20.6	
12	CBN, 2020, [225]	Market-1501	94.3/83.6	Yes
		DukeMTMC-ReId	84.8/70.1	
13	CD-ReID, 2020, [226]	Market-1501	88.1/71.5	Yes
		DukeMTMC-ReId	79.5/65.2	

14	DG-Net++, [227]	2020,	Market-1501	83.1/64.6	Yes
			DukeMTMC-ReId	78.9/63.8	
			MSMT-17	48.8/22.1	
15	TLift, 2020, [268]		Market-1501	88.4/76.0	Yes
			DukeMTMC-ReId	82.2/78.4	
16	D-MMD, 2020, [250]		Market-1501	72.8/50.8	Yes
			DukeMTMC-ReId	68.8/51.6	
			MSMT-17	34.4/15.3	
17	SSDAL, 2016, [245]		Market-1501	49.0/25.8	No
			VIPeR	43.5/–	
			PRID-2011	22.6/–	
18	MARS, 2016, [17]		MARS	68.3/49.3	No
			iLIDS	53.0/–	
			PRID-2011	77.3/–	
19	DMG-Net, [228]	2021,	Person-30K	84.23/72.19	No
20	RDSBN, 2021, [251]		Market-1501	94.8/86.0	No
			DukeMTMC-ReId	82.1/68.9	
			MSMT-17	64.7/34.9	
21	MetaBIN, 2021, [229]		Market-1501	69.2/35.9	yes
			DukeMTMC-ReId	55.2/33.1	
			VIPeR	59.9/68.6	
			iLIDS	81.3/87.0	
			PRID-2011	81.0/72.4	
22	RaMoE, 2021, [230]		Market-1501	82.0/56.5	No
			DukeMTMC-ReId	73.6/56.9	
			CUHK-03	36.6/35.5	
			VIPeR	63.4/72.2	
			MSMT-17	34.1/13.5	
			iLIDS	88.4/92.3	
23	LUPerson, [252]	2021,	Market-1501	97.0/92.0	No
			DukeMTMC-ReId	91.9/84.1	
			CUHK-03	81.9/79.6	
			MSMT-17	86.6/68.8	
24	LReID, 2021, [272]		Market-1501	87.0/74.8	Yes
			DukeMTMC-ReId	80.1/68.3	
			CUHK-03	56.6/50.8	
			MSMT-17	54.1/27.9	
25	IICS, 2021, [253]		Market-1501	89.5/72.9	Yes
			DukeMTMC-ReId	80.0/64.4	
			MSMT-17	56.4/26.9	
26	DSCE, 2021, [254]		Market-1501	83.9/61.7	Yes
			DukeMTMC-ReId	73.8/53.8	
			MSMT-17	35.2/15.5	
27	ADC-2OIB, [231]	2021,	Market-1501	94.8/87.7	No
			DukeMTMC-ReId	87.4/74.9	
			CUHK-03(Labeled)	80.6/79.3	
			CUHK-03(Detected)	81.3/84.1	

28	RLCC, 2021, [255]	Market-1501 DukeMTMC-ReId MSMT-17	90.8/77.7 83.2/69.2 56.5/27.9	No
29	UnrealPerson, 2021, [232]	Market-1501 DukeMTMC-ReId MSMT-17	93.0/80.2 88.3/75.2 68.2/34.8	Yes
30	M3L, 2021, [256]	Market-1501 DukeMTMC-ReId CUHK-03 MSMT-17	75.9/50.2 69.2/51.1 33.1/32.1 36.9/14.7	Yes
31	GLT, 2021, [257]	Market-1501 DukeMTMC-ReId MSMT-17	92.2/79.5 82.0/69.2 59.5/27.7	No
32	SNR, 2020, [221]	Market-1501 DukeMTMC-ReId	85.5/65.9 78.2/61.6	No
33	AD-Cluster, 2020, [267]	Market-1501 DukeMTMC-ReId	86.7/68.3 72.6/54.1	No
34	ATNet, 2019, [73]	Market-1501 DukeMTMC-ReId	45.1/24.9 55.7/25.6	No
35	AANet, 2019, [275]	Market-1501 DukeMTMC-ReId	95.1/92.38 90.36/36.87	No
36	PatchNet, 2019, [243]	Market-1501 DukeMTMC-ReId	68.5/40.1 72.0/53.2	Yes
37	EANet, 2019, [153]	Market-1501 DukeMTMC-ReId CUHK-03	94.5/85.6 87.5/74.6 72.5/66.8	Yes
38	DG-Net, 2019, [219]	Market-1501 DukeMTMC-ReId MSMT-17	94.8/86.0 86.6/74.8 77.2/52.3	No
39	CV-MIML, 2019, [266]	DukeMTMC-VideoReId MARS iLIDS PRID-2011	78.05/59.53 66.88/55.16 60.0/56.01 72.0/70.78	No
40	SPGAN, 2018, [271]	Market-1501 DukeMTMC-ReId	58.1/26.9 46.9/26.4	No
41	DuATM, 2018, [276]	Market-1501 DukeMTMC-ReId MARS	91.42/76.62 81.82/64.58 78.74/62.26	No
42	QDNet, 2017, [222]	CUHK-01 CUHK-03 VIPeR	81.0/– 75.53/– 49.05/–	No
43	DGD, 2016, [223]	CUHK-01 CUHK-03 VIPeR iLIDS PRID-2011	66.6/– 75.3/– 38.6/– 64.6/– 64.0/–	Yes
44	UMDL, 2016, [247]	CUHK-01 VIPeR iLIDS PRID-2011	27.1/– 31.5/– 49.3/– 24.2/–	Yes

45	SAL, 2015, [246]	CUHK-01 CUHK-03	22.4/– 29.3/–	No
46	4S-Net, 2020, [234]	Market-1501 DukeMTMC-ReId VIPeR iLIDS	91.6/75.7 82.4/77.3 71.4 /– 84.8/–	No
47	PL-Net, 2019, [244]	Market-1501 CUHK-03 VIPeR	88.2/69.3 82.75/– 56.65/–	No
48	kLDFA, 2016, [182]	CUHK-01 VIPeR	57.28/– 38.37/–	No
49	CGAN-TM, 2020, [270]	Market-1501 DukeMTMC-ReId	64.43/31.33 54.85/32.85	No
50	GPP-ReID, 2021, [236]	Market-1501 DukeMTMC-ReId MSMT-17	90.6/78.6 81.3/67.9 53.5/24.6	No
51	CI-CNN, 2020, [237]	Market-1501 DukeMTMC-ReId MARS	94.26/89.54 87.6/81.3 87.3/78.8	No
52	DRM, 2021, [238]	Market-1501 DukeMTMC-ReId MSMT-17	90.9/78.0 82.1/67.7 55.0/26.9	No
53	HCC-GCNs, 2021, [260]	Market-1501 DukeMTMC-ReId MSMT-17	91.2/78.9 81.2/67.5 57.4/28.4	No
54	IIA, 2020, [279]	Market-1501 DukeMTMC-ReId CUHK-03 (Labeled) CUHK-03 (Detected)	98.2/96.0 94.4/91.8 84.93/86.58 80.14/82.72	No
55	VOLTA, 2020, [261]	MARS DukeMTMC-VideoReId	66.7/51.9 89.2/85.9	No
56	MMFA-AAE, 2021, [239]	VIPeR MSMT-17 iLIDS	58.4/– 46.0/20.7 84.8/–	No
57	SAL, 2020, [262]	Market-1501 DukeMTMC-ReId CUHK-03 (Labeled) CUHK-03 (Detected)	68.7/41.9 70.8/51.4 31.0/35.7 28.7/34.4	No
58	PREST, 2021, [172]	Market-1501 DukeMTMC-ReId	82.5/62.4 74.4/56.1	No
59	MLOL, 2021, [263]	Market-1501 DukeMTMC-ReId MSMT-17	86.6/70.9 83.1/69.8 48.3/22.4	No
60	CCAN, 2020, [277]	Market-1501 DukeMTMC-ReId CUHK-03 (Labeled) CUHK-03 (Detected) MSMT-17	94.6/87.0 87.2/76.8 75.2/72.9 73.0/70.7 76.3/53.6	Yes

61	JPIL, 2020, [259]	Market-1501 DukeMTMC-ReId	73.5/48.2 74.7/55.8	No
62	CDL, 2018, [131]	CUHK-01 VIPeR iLIDS	78.17/– 66.39/– 45.1/–	Yes
63	DECAMEL, 2020, [202]	Market-1501 CUHK-01 CUHK-03 MSMT-17	60.24/32.44 65.81/– 38.27/– 90.34/11.13	No
64	UTAL, 2020, [264]	Market-1501 DukeMTMC-ReId CUHK-03 MSMT-17 MARS iLIDS	69.2/46.2 62.3/44.6 56.3/42.3 31.4/13.1 49.9/35.2 35.1/–	No
65	HAN, 2019, [278]	Market-1501 DukeMTMC-ReId CUHK-03 (Labeled) CUHK-03 (Detected) MSMT-17	94.2/83.4 80.6/64.1 46.5/46.1 47.5/45.5 60.1/32.6	No
66	SBSGAN, 2021, [120]	Market-1501 DukeMTMC-ReId	87.9/80.0 79.7/71.5	No
67	3D-GAT, 2021, [240]	Market-1501 DukeMTMC-ReId	94.1/81.5 85.5/71.2	No
68	pLMNN, 2018, [241]	CUHK-01 VIPeR	53.5/– 46.5/–	No
69	SADL, 2020, [242]	Market-1501 DukeMTMC-ReId CUHK-03 MSMT-17	60.7/26.6 50.2/28.1 75.42/– 30.6/–	No
70	UDAM, 2020, [273]	Market-1501 DukeMTMC-ReId	75.8/53.7 68.4/49.0	Yes
71	LADL, 2020, [265]	CUHK-01	57.67/–	No
72	IVD-ReID, 2019, [198]	MARS iLIDS	48.0/– 65.0/–	No

Table 13 comprehensively shows the SOTA results achieved on each dataset.

6. Discussion & Future Trends

In this systematic review, 230+ articles are reviewed that were published from January 2015 till October 2021 focused on the challenges faced for person re-id. In all these papers, the specified challenges have been addressed and the achieved state-of-the-results have been summarized in this review article. We have grouped the articles into few categories and have critically analysed their impact on the obtained results. The limitations along with the datasets used in the published articles have also been reviewed against each challenge.

6.1. Impact of Automated Person Re-id on Society

Generally, due to the scarcity of security resources as well as the lack of technological advancements in third world countries, the traditional law and order system does not meet the needs of the public, hence failing to build the people's confidence. The automated surveillance systems *i.e.* person re-id aims to improve the quality of the lives of common

Table 13

SOTA results obtained on each challenge against each dataset.

Sr.No	Dataset	SOTA R1-results	Paper cited	Challenge addressed	Venue
1	Market-1501	98.3	[119]	Background	TCSV-2020
2	DukeMTMC-ReID	94.7	[119]	Background	TCSV-2020
3	CUHK01	98.73	[165]	Scale	PAMI-2020
4	CUHK03	97.3	[91]	Pose	PAMI-2021
5	VIPeR	71.4	[234]	Generalization	CVIU-2020
6	MSMT-17	87.7	[141]	Misalignment	PR-2021
7	PRID-2011	95.9	[99]	Pose	PR-2021
8	MARS	91.5	[67]	Occlusion	ICCV-2021
9	DukeMTMC-Video ReID	98.3	[67]	Occlusion	ICCV-2021
10	iLIDS	92	[142](arxiv)	Misalignment	ICCV-2021

people by providing a sustainable living environment to them. Since the assurance of security and implementation of law and order are from the basic needs of human beings, the person re-id solutions can assist the law enforcement authorities in providing enhanced preemptive security and quick implementation of the law and order. Moreover, in case of any adverse happening, the re-id solutions can greatly assist the security officials in rapid response and quick resolution of the security issues and can prevent the delays caused by manual video forensic analytic. In order to trace people in a camera network, faces can only be used for recognition only if the subject is close enough and facing towards the camera. But usually in CCTV footage this is not the case, people are captured at variant poses and viewpoints where their faces are not clear. Therefore, the features found in a person's entire body (like clothing, height etc) are more useful to identify a person across different cameras of the network.

6.2. Deep Learning Conjecture

Since the evolution of deep learning methods in 2015, the computer vision research community immediately shifted from hand-crafted machine learning research to the deep learning based algorithms. Soon after the availability of a very large scale vision based dataset *i.e.* ImageNet and its pre-trained weights, like many other research domain, deep person re-id solutions came into existence. Meanwhile, the medium to large scale re-id benchmarks were proposed so that the specialized custom re-id solutions could be developed. Beginning from the transfer learning methods using the generic deep architectures like AlexNet, ResNet *etc.*, the re-id research quickly got independence in designing much sophisticated solutions due to the availability of large scale person re-id benchmarks. For many years, the convolution neural networks served as a strong backbone for re-id solutions. The CNNs are used to learn a variety of person representations, *i.e.* the global person representations, local parts based person representations, semantics based person representations, attributes driven representations *etc.* The CNN architectures with single stream/ branch was common in the start, however with the passage of time, multi-streamed architectures are proposed for person re-id, where each stream targets a different perspective.

Later on, with the development of attention based mechanism for vision problems, the same were extensively explored to develop the re-id solutions. Most of the attention based re-id solutions are developed on the backbone of CNNs and are multi-stream architectures to capture the various types of attention features *i.e.* spatial attention, temporal attention, channel wise attention *etc.* The attention based re-id solutions performed really well for various re-id challenges in comparison with the methods that do not involve computation of the attention.

Since the deep-learning is still evolving and new SOTA backbone architectures are being developed by the research community with every passing year, it drives the whole research community to new dimensions. Transformers and its variants are SOTA for the language problems since long but due to certain limitations, these were not used for vision problems in a holistic way. Recently, the development of vision transformer (in year 2021) was a great break-through in the vision research, and opened up the new ways in vision research. The vision transformer outperforms the counterpart

CNN baseline deep architectures with a great margin for various vision problems. Since then, the re-id solutions with transformers backbone are gaining popularity among the research community.

This comprehensive study on person re-identification research unfolds few interesting aspects for the future re-id researches. While exploring the re-id solutions for the foremost common re-id challenge, i.e., the occlusion, it has been observed that the top three re-id solutions [67], [70], [71] that optimally re-identify the occluded persons are all attention based approaches. These attention based mechanisms are strengthened by the use of 3D convolutional architecture, pyramid architecture and the memory units. The different perspectives of the spatial and temporal features are learnt to capture the the dynamic and static information of a person. The temporal features compensate for the occluded spatial regions and enhance the performances of the re-id solutions.

The pose and view point variations make the person re-id quite challenging especially in case of inter-class differences (where people of different classes appear similar under different camera acquisitions) due to the similar appearances. This study highlights that for this particular challenge, both the CNN and attention based approaches are comparable. For instance, among the top three re-id solutions [90], [89] and [107], only [90] presented the approach which is based on learning attention weights while [89] employed the higher level semantic information to generate multi-level feature maps and [107] presented a novel Kronecker product matching operation to perform the re-identification. Generally, these approaches work to handle the misalignment by learning aligned and misaligned parts of person images. Local parts based dynamic feature learning addresses the misalignment issues either by focusing on motion specific and joint specific information in person images.

The cluttered background, if not efficiently removed or suppressed, inversely effects the performance of re-id solution. Since the attention based mechanism inherently focus and highlight the attentive parts/ regions on a person image, the top three re-id solutions which perform the best even in case of cluttered background also employed attention based mechanisms to exclude the background information [115] and learns multi-level attention in different ways [118] and fuse the information learnt from various levels and branches [119]. Generally, these approaches capture the foreground attention features using either the encoder decoder architecture or the multi-level attention modules. Learning the binary masks is common practice to exclude the cluttered background from the person foreground. However, these approaches do not aid the scenario where the background information might play vital role in the identification of a person.

The orientation of cameras and different viewing angles result in the misalignment of person images. To tackle Either the re-id solutions perform images alignment by proposing part-based re-id solutions or align the image regions by proposing various sophisticated algorithms. Top three approaches to address this issue are both the CNN based and the attention based. [142], an encoder-decoder based architectures, is a hybrid approach that takes the advantage of both CNN and attention based mechanism to handle the misalignment efficiently and outperformed all other solutions for video-based re-id benchmarks. [47] is an attention mechanism that handled the misalignment by learning the channel wise and position/spatial information and [141] handled the misalignment by part-guided graph convolution network. Generally, the local parts/ patches based learning in addition to global feature learning aids to address the misalignment issues in person re-id.

The differences in scales of captured images is generally seen in CCTV footage due to the variations in the distances of cameras from the targets. The scale differences are handled by various re-id solutions, however, the attention based approaches outperformed the rest of re-id solutions. The top three solutions either used the multi-scale attention pyramid [165] or divided the image into multiple local parts and then learnt the attention [169], or [171] extracted the holistic and local feature maps using multi-scale omni-bearing attention network. The re-id solutions that support multi-scale re-id learn the person features at multiple scales through multi-sized convolutional layers or branches. And then aggregate the learned information into final person descriptor. For video based benchmarks, the temporal information alleviate the multi-scale re-id solutions.

An effective re-id solution handles the illumination variations efficiently such that the changes in the color of a person's dress due to variant lightening conditions may not result in false re-identification of an entity. Similarly the view-point variations need significant attention by the underlying re-id solution. The best performance is attained by an attention based mechanism [203] to address these re-id challenges. However, few other convolutions based architectures [196] and [185] performed well to address the view-point variations.

Since the surveillance cameras work 24/7 and capture the person images from a distance, generally the low image resolution is observed in CCTV footage. This results in another re-id challenge i.e. to identify a person correctly in low resolution images. Recently, [213] proposed a resolution-aware re-id framework that works very well and follows the teacher-student learning mechanism.

Lastly, the cross-domain person re-identification is quite challenging with a huge room for improvement. Due to numerous re-id challenges, already discussed in previous sections, it is difficult to re-identify the people of totally different camera network. A well generalized re-id solution is highly desirable for cross-domain re-identification. A lot of research has been carried out to solve this problem and recently a generalized re-id solution [252] used the temperature features to propose a contrastive learning framework. Another attention aggregation formulation was designed by [279] for a generalized re-id solution.

6.3. Concluding Remarks

In general, we can emphasize that attention based re-id solutions are gaining more interest in the research community with their promising performances. Also, since the self-attention based vision research is just in its beginning phase, the full strength of these approaches are yet to be researched and analyzed.

During the recent years, the re-id research is at its best for few initially proposed re-id datasets *i.e.* Market1501, DukeMTMC-Reid etc, which were captured from the public places with controlled environment, hence do not depict the real world scenario. The customized re-id algorithms addressed most of the re-id challenges effectively due to the medium level of complexity for these benchmarks. Among various re-id challenges, the pose variations remained most popular among research community, as it is the most common re-id issue with significant impact on the performance of re-id solutions. The newly proposed re-id benchmarks (*i.e.* MSMT17 etc), captured from the complex scenes with a large number of indoor and outdoor cameras and closer to the real world complex scenarios, therefore these need more sophisticated re-id solutions to solve the real world problems.

Since the start of the re-id research journey, the majority of re-id solutions work for re-id of people from a single benchmark for unseen identities. However, far less research is done to propose the re-id solution for cross-domain benchmarks. The major reason is that, the re-id research is still evolving and needs sophisticated research solutions even for the independent re-id benchmarks. For the recently proposed complex and large-scale re-id benchmarks, the performance of the re-id solutions need significant improvement.

Declarations

The authors declare no conflict of interest.

References

- [1] Zheng Wang, Ruimin Hu, Chao Liang, Yi Yu, Junjun Jiang, Mang Ye, Jun Chen, and Qingming Leng. Zero-shot person re-identification via cross-view consistency. *IEEE Transactions on Multimedia*, 18(2):260–272, 2015.
- [2] Yujiang Wang, Jie Shen, Stavros Petridis, and Maja Pantic. A real-time and unsupervised face re-identification system for human-robot interaction. *Pattern Recognition Letters*, 128:559–568, 2019.
- [3] Hanxiao Wang, Shaogang Gong, Xiatian Zhu, and Tao Xiang. Human-in-the-loop person re-identification. In *European conference on computer vision*, pages 405–422. Springer, 2016.
- [4] Liang Zheng, Yi Yang, and Alexander G Hauptmann. Person re-identification: Past, present and future. *arXiv preprint arXiv:1610.02984*, 2016.
- [5] Kejun Wang, Haolin Wang, Meichen Liu, Xianglei Xing, and Tian Han. Survey on person re-identification based on deep learning. *CAAI Transactions on Intelligence Technology*, 3(4):219–227, 2018.
- [6] Di Wu, Si-Jia Zheng, Xiao-Ping Zhang, Chang-An Yuan, Fei Cheng, Yang Zhao, Yong-Jun Lin, Zhong-Qiu Zhao, Yong-Li Jiang, and De-Shuang Huang. Deep learning-based methods for person re-identification: A comprehensive review. *Neurocomputing*, 337:354–371, 2019.
- [7] Mang Ye, Jianbing Shen, Gaojie Lin, Tao Xiang, Ling Shao, and Steven CH Hoi. Deep learning for person re-identification: A survey and outlook. *arXiv preprint arXiv:2001.04193*, 2020.
- [8] Qingming Leng, Mang Ye, and Qi Tian. A survey of open-world person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(4):1092–1108, 2019.
- [9] University of Oxford Ottawa Hospital Research Institute. Prisma: Transparent reporting of systematic reviews and meta-analyses.
- [10] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*, pages 1116–1124, 2015.
- [11] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In *European conference on computer vision*, pages 17–35. Springer, 2016.
- [12] Wei Li, Rui Zhao, and Xiaogang Wang. Human reidentification with transferred metric learning. In *Asian conference on computer vision*, pages 31–44. Springer, 2012.
- [13] Wei Li, Rui Zhao, Tong Xiao, and Xiaogang Wang. Deepreid: Deep filter pairing neural network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 152–159, 2014.

- [14] Douglas Gray, Shane Brennan, and Hai Tao. Evaluating appearance models for recognition, reacquisition, and tracking. In *Proc. IEEE international workshop on performance evaluation for tracking and surveillance (PETS)*, volume 3, pages 1–7. Citeseer, 2007.
- [15] Martin Hirzer, Csaba Beleznaï, Peter M Roth, and Horst Bischof. Person re-identification by descriptive and discriminative classification. In *Scandinavian conference on Image analysis*, pages 91–102. Springer, 2011.
- [16] Taiqing Wang, Shaogang Gong, Xiatian Zhu, and Shengjin Wang. Person re-identification by video ranking. In *European conference on computer vision*, pages 688–703. Springer, 2014.
- [17] Liang Zheng, Zhi Bie, Yifan Sun, Jingdong Wang, Chi Su, Shengjin Wang, and Qi Tian. Mars: A video benchmark for large-scale person re-identification. In *European Conference on Computer Vision*, pages 868–884. Springer, 2016.
- [18] Longhui Wei, Shiliang Zhang, Wen Gao, and Qi Tian. Person transfer gan to bridge domain gap for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 79–88, 2018.
- [19] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2009.
- [20] Afshin Dehghan, Shayan Modiri Assari, and Mubarak Shah. Gmmcp tracker: Globally optimal generalized maximum multi clique problem for multiple object tracking. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4091–4099, 2015.
- [21] Zhedong Zheng, Liang Zheng, and Yi Yang. A discriminatively learned cnn embedding for person reidentification. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 14(1):1–20, 2017.
- [22] Yutian Lin, Liang Zheng, Zhedong Zheng, Yu Wu, Zhilan Hu, Chenggang Yan, and Yi Yang. Improving person re-identification by attribute and identity learning. *Pattern Recognition*, 95:151–161, 2019.
- [23] Zhangping He, Cheolkon Jung, Qingtao Fu, and Zhenhong Zhang. Deep feature embedding learning for person re-identification based on lifted structured loss. *Multimedia Tools and Applications*, 78(5):5863–5880, Mar 2019.
- [24] Jingjing Wu, Jianguo Jiang, Meibin Qi, and Hao Liu. Independent metric learning with aligned multi-part features for video-based person re-identification. *Multimedia Tools and Applications*, 78(20):29323–29341, Oct 2019.
- [25] Alexander Hermans, Lucas Beyer, and Bastian Leibe. In defense of the triplet loss for person re-identification. *arXiv preprint arXiv:1703.07737*, 2017.
- [26] Xinyu Ou, Qianzhi Ma, and Yijin Wang. Improving person re-identification by multi-task learning. *Multimedia Tools and Applications*, 78(19):28257–28283, Oct 2019.
- [27] Tiezhu Li, Lijuan Sun, Chong Han, and Jian Guo. Person re-identification using salient region matching game. *Multimedia Tools and Applications*, 77(16):21393–21415, Aug 2018.
- [28] N. Pervaiz, M.M. Fraz, and M. Shahzad. Hierarchical refined local associations for robust person re-identification. In *IEEE International Conference on Robotics and Automation, Islamabad, Pakistan*, volume 3, Nov 2019.
- [29] Saadia Batool, Muhammad Zeeshan Ali, Muhammad Shahzad, and Muhammad Moazam Fraz. End to end person re-identification for automated visual surveillance. In *IEEE International Conference on Image Processing, Applications and Systems, IPAS 2018, Sophia Antipolis, France, December 12-14, 2018*, pages 220–225, 2018.
- [30] Wajeeha Ansar, M. M. Fraz, M. Shahzad, I. Gohar, Sajid Javed, and Soon Ki Jung. Two stream deep CNN-RNN attentive pooling architecture for video-based person re-identification. In *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications - 23rd Iberoamerican Congress, CIARP 2018, Madrid, Spain, November 19-22, 2018, Proceedings*, pages 654–661, 2018.
- [31] Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Cross-view asymmetric metric learning for unsupervised person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 994–1002, 2017.
- [32] Douglas Gray and Hai Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *European conference on computer vision*, pages 262–275. Springer, 2008.
- [33] Xiaojuan Wang, Wei-Shi Zheng, Xiang Li, and Jianguo Zhang. Cross-scenario transfer person reidentification. *IEEE Transactions on Circuits and Systems for Video Technology*, 26(8):1447–1460, 2015.
- [34] Michela Farenzena, Loris Bazzani, Alessandro Perina, Vittorio Murino, and Marco Cristani. Person re-identification by symmetry-driven accumulation of local features. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 2360–2367. IEEE, 2010.
- [35] Elyor Kodirov, Tao Xiang, Zhenyong Fu, and Shaogang Gong. Person re-identification by unsupervised graph learning. In *European conference on computer vision*, pages 178–195. Springer, 2016.
- [36] Zhiyuan Shi, Timothy M Hospedales, and Tao Xiang. Transferring a semantic representation for person re-identification and search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4184–4193, 2015.
- [37] Saba Mumtaz, Naima Mubariz, Shahzad Saleem, and M. M. Fraz. Weighted hybrid features for person re-identification. In *Seventh International Conference on Image Processing Theory, Tools and Applications, IPTA 2017, Montreal, QC, Canada, November 28 - December 1, 2017*, pages 1–6, 2017.
- [38] Naima Mubariz, Saba Mumtaz, Mian M. Hamayun, and M. M. Fraz. Optimization of person re-identification through visual descriptors. In *Proceedings of (VISIGRAPP 2018) - Volume 4: VISAPP, Funchal, Madeira, Portugal, January 27-29, 2018.*, pages 348–355, 2018.
- [39] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30:5998–6008, 2017.
- [40] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [42] Nazia Perwaiz, Muhammad Moazam Fraz, and Muhammad Shahzad. Person re-identification using hybrid representation reinforced by metric learning. *IEEE Access*, 6:77334–77349, 2018.

- [43] Rao Faizan, Muhammad Moazam Fraz, and Muhammad Shahzad. Iab-net: Informative and attention based person re-identification. In *2021 International Conference on Digital Futures and Transformative Technologies (ICoDT2)*, pages 1–5. IEEE, 2021.
- [44] Nazia Perwaiz, Muhammad Moazam Fraz, and Muhammad Shahzad. Stochastic attentions and context learning for person re-identification. *PeerJ Computer Science*, 7:e447, 2021.
- [45] N Perwaiz, MM Fraz, and M Shahzad. Smart visual surveillance: Proactive person re-identification instead of impulsive person search. In *2020 IEEE 23rd International Multitopic Conference (INMIC)*, pages 1–6. IEEE, 2020.
- [46] Guangcong Wang, Jianhuang Lai, Peigen Huang, and Xiaohua Xie. Spatial-temporal person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8933–8940, 2019.
- [47] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. Abd-net: Attentive but diverse person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8351–8361, 2019.
- [48] Guangyi Chen, Chunze Lin, Liangliang Ren, Jiwen Lu, and Jie Zhou. Self-critical attention learning for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9637–9646, 2019.
- [49] Yifan Chen, Han Wang, Xiaolu Sun, Bin Fan, and Chu Tang. Deep attention aware feature learning for person re-identification. *arXiv preprint arXiv:2003.00517*, 2020.
- [50] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2285–2294, 2018.
- [51] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, and Sylvain Gelly. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [52] Hao Luo, Wei Jiang, Xing Fan, and Chi Zhang. Stnreid: Deep convolutional networks with pairwise spatial transformer networks for partial person re-identification. *IEEE Transactions on Multimedia*, 22(11):2905–2913, 2020.
- [53] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Learning deep context-aware features over body and latent parts for person re-identification. *CoRR*, abs/1710.06555, 2017.
- [54] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 480–496, 2018.
- [55] Wei-Shi Zheng, Xiang Li, Tao Xiang, Shengcai Liao, Jianhuang Lai, and Shaogang Gong. Partial person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4678–4686, 2015.
- [56] Lingxiao He, Yinggang Wang, Wu Liu, He Zhao, Zhenan Sun, and Jiashi Feng. Foreground-aware pyramid reconstruction for alignment-free occluded person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8450–8459, 2019.
- [57] Lingxiao He and Wu Liu. Guided saliency feature learning for person re-identification in crowded scenes. In *European Conference on Computer Vision*, pages 357–373. Springer, 2020.
- [58] Guan'an Wang, Shuo Yang, Huanan Liu, Zhicheng Wang, Yang Yang, Shuliang Wang, Gang Yu, Erjin Zhou, and Jian Sun. High-order information matters: Learning relation and topology for occluded person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6449–6458, 2020.
- [59] Cheng Yan, Guansong Pang, Jile Jiao, Xiao Bai, Xuetao Feng, and Chunhua Shen. Occluded person re-identification with single-scale global representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11875–11884, 2021.
- [60] Jinrui Yang, Jiawei Zhang, Fufu Yu, Xinyang Jiang, Mengdan Zhang, Xing Sun, Ying-Cong Chen, and Wei-Shi Zheng. Learning to know where to see: A visibility-aware approach for occluded person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11885–11894, 2021.
- [61] Cairong Zhao, Xinbi Lv, Shuguang Dou, Shanshan Zhang, Jun Wu, and Liang Wang. Incremental generative occlusion adversarial suppression network for person reid. *IEEE Transactions on Image Processing*, 30:4212–4224, 2021.
- [62] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Feature completion for occluded person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [63] Xiaokang Zhang, Yan Yan, Jing-Hao Xue, Yang Hua, and Hanzi Wang. Semantic-aware occlusion-robust network for occluded person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [64] Jiaxu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. Pose-guided feature alignment for occluded person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 542–551, 2019.
- [65] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Vrstc: Occlusion-free video person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7183–7192, 2019.
- [66] Shizhen Zhao, Changxin Gao, Jun Zhang, Hao Cheng, Chuchu Han, Xinyang Jiang, Xiaowei Guo, Wei-Shi Zheng, Nong Sang, and Xing Sun. Do not disturb me: Person re-identification under the interference of other pedestrians. In *European Conference on Computer Vision*, pages 647–663. Springer, 2020.
- [67] Yingquan Wang, Pingping Zhang, Shang Gao, Xia Geng, Hu Lu, and Dong Wang. Pyramid spatial-temporal aggregation for video-based person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12026–12035, 2021.
- [68] Ruimao Zhang, Jingyu Li, Hongbin Sun, Yuying Ge, Ping Luo, Xiaogang Wang, and Liang Lin. Scan: Self-and-collaborative attention network for video person re-identification. *IEEE Transactions on Image Processing*, 28(10):4870–4882, 2019.
- [69] Guangyi Chen, Jiwen Lu, Ming Yang, and Jie Zhou. Spatial-temporal attention-aware learning for video-based person re-identification. *IEEE Transactions on Image Processing*, 28(9):4192–4205, 2019.
- [70] Abhishek Aich, Meng Zheng, Srikrishna Karanam, Terrence Chen, Amit K Roy-Chowdhury, and Ziyang Wu. Spatio-temporal representation factorization for video-based person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 152–162, 2021.
- [71] Chanho Eom, Geon Lee, Junghyup Lee, and Bumsub Ham. Video-based person re-identification with spatial and temporal memory networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 12036–12045, 2021.

- [72] Yulin Li, Jianfeng He, Tianzhu Zhang, Xiang Liu, Yongdong Zhang, and Feng Wu. Diverse part discovery: Occluded person re-identification with part-aware transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2898–2907, 2021.
- [73] Jiawei Liu, Zheng-Jun Zha, Di Chen, Richang Hong, and Meng Wang. Adaptive transfer network for cross-domain person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7202–7211, 2019.
- [74] Ruibing Hou, Bingpeng Ma, Hong Chang, Xinqian Gu, Shiguang Shan, and Xilin Chen. Interaction-and-aggregation network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9317–9326, 2019.
- [75] George Papandreou, Tyler Zhu, Liang-Chieh Chen, Spyros Gidaris, Jonathan Tompson, and Kevin Murphy. Personlab: Person pose estimation and instance segmentation with a bottom-up, part-based, geometric embedding model. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 269–286, 2018.
- [76] Chi Su, Jianing Li, Shiliang Zhang, Junliang Xing, Wen Gao, and Qi Tian. Pose-driven deep convolutional model for person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 3960–3969, 2017.
- [77] Muhammed Kocabas, Salih Karagoz, and Emre Akbas. Multiposenet: Fast multi-person pose estimation using pose residual network. In *Proceedings of the European conference on computer vision (ECCV)*, pages 417–433, 2018.
- [78] Yeong-Jun Cho and Kuk-Jin Yoon. Improving person re-identification via pose-aware multi-shot matching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1354–1362, 2016.
- [79] Rodolfo Quispe and Helio Pedrini. Improved person re-identification based on saliency and semantic parsing with deep neural network models. *Image and Vision Computing*, 92:103809, 2019.
- [80] Yiqiang Chen, Stefan Duffner, Andrei Stoian, Jean-Yves Dufour, and Atila Baskurt. Deep and low-level feature based attribute learning for person re-identification. *Image and Vision Computing*, 79:25–34, 2018.
- [81] Niall McLaughlin, Jesus Martinez del Rincon, and Paul C Miller. Person reidentification using deep convnets with multitask learning. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(3):525–539, 2016.
- [82] Jin Wang, Zheng Wang, Changxin Gao, Nong Sang, and Rui Huang. Deeplist: Learning deep features with adaptive listwise constraint for person reidentification. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(3):513–524, 2016.
- [83] Young-Gun Lee, Shen-Chi Chen, Jenq-Neng Hwang, and Yi-Ping Hung. An ensemble of invariant features for person reidentification. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(3):470–483, 2016.
- [84] Jianing Li, Shiliang Zhang, Qi Tian, Meng Wang, and Wen Gao. Pose-guided representation learning for person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [85] Shoubiao Tan, Feng Zheng, Li Liu, Jungong Han, and Ling Shao. Dense invariant feature-based support vector ranking for cross-camera person reidentification. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(2):356–363, 2016.
- [86] Igor Barros Barbosa, Marco Cristani, Barbara Caputo, Aleksander Rognhaugen, and Theoharis Theoharis. Looking beyond appearances: Synthetic training data for deep cnns in re-identification. *Computer Vision and Image Understanding*, 167:50–62, 2018.
- [87] Zhiyong Li, Jingyi Lv, Ying Chen, and Jin Yuan. Person re-identification with part prediction alignment. *Computer Vision and Image Understanding*, 205:103172, 2021.
- [88] Jiahang Yin, Ancong Wu, and Wei-Shi Zheng. Fine-grained person re-identification. *International journal of computer vision*, 128(6):1654–1672, 2020.
- [89] Qinqin Zhou, Bineng Zhong, Xiangyuan Lan, Gan Sun, Yulun Zhang, Baochang Zhang, and Rongrong Ji. Fine-grained spatial alignment model for person re-identification with focal triplet loss. *IEEE Transactions on Image Processing*, 29:7578–7589, 2020.
- [90] Pingyu Wang, Zhicheng Zhao, Fei Su, Xingyu Zu, and Nikolaos V Boulgouris. Horeid: Deep high-order mapping enhances pose alignment for person re-identification. *IEEE Transactions on Image Processing*, 30:2908–2922, 2021.
- [91] Yantao Shen, Tong Xiao, Shuai Yi, Dapeng Chen, Xiaogang Wang, and Hongsheng Li. Person re-identification with deep kronecker-product matching and group-shuffling random walk. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1649–1665, 2021.
- [92] Chi Su, Shiliang Zhang, Junliang Xing, Wen Gao, and Qi Tian. Multi-type attributes driven multi-camera person re-identification. *Pattern Recognition*, 75:77–89, 2018.
- [93] M Saquib Sarfraz, Arne Schumann, Andreas Eberle, and Rainer Stiefelhausen. 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*, pages 420–429, 2018.
- [94] Liang Zheng, Yujia Huang, Huchuan Lu, and Yi Yang. Pose-invariant embedding for deep person re-identification. *IEEE Transactions on Image Processing*, 28(9):4500–4509, 2019.
- [95] Shih-En Wei, Varun Ramakrishna, Takeo Kanade, and Yaser Sheikh. Convolutional pose machines. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 4724–4732, 2016.
- [96] Jinxian Liu, Bingbing Ni, Yichao Yan, Peng Zhou, Shuo Cheng, and Jianguo Hu. Pose transferrable person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4099–4108, 2018.
- [97] Xuecheng Nie, Jiashi Feng, Junliang Xing, and Shuicheng Yan. Pose partition networks for multi-person pose estimation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 684–699, 2018.
- [98] Yiming Wu, Omar El Farouk Bourahla, Xi Li, Fei Wu, Qi Tian, and Xue Zhou. Adaptive graph representation learning for video person re-identification. *IEEE Transactions on Image Processing*, 29:8821–8830, 2020.
- [99] Xiaoqiang Hu, Dan Wei, Ziyang Wang, Jianglin Shen, and Hongjuan Ren. Hypergraph video pedestrian re-identification based on posture structure relationship and action constraints. *Pattern Recognition*, 111:107688, 2021.
- [100] Dapeng Chen, Zejian Yuan, Badong Chen, and Nanning Zheng. Similarity learning with spatial constraints for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1268–1277, 2016.
- [101] Dapeng Chen, Zejian Yuan, Gang Hua, Nanning Zheng, and Jingdong Wang. Similarity learning on an explicit polynomial kernel feature map for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1565–1573,

2015.

- [102] Le An, Mehran Kafai, Songfan Yang, and Bir Bhanu. Person reidentification with reference descriptor. *IEEE Transactions on Circuits and Systems for Video Technology*, 26(4):776–787, 2015.
- [103] Yifan Sun, Liang Zheng, Yali Li, Yi Yang, Qi Tian, and Shengjin Wang. Learning part-based convolutional features for person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 2019.
- [104] Yeong-Jun Cho and Kuk-Jin Yoon. Pamm: Pose-aware multi-shot matching for improving person re-identification. *IEEE Transactions on Image Processing*, 27(8):3739–3752, 2018.
- [105] Wenjie Yang, Houjing Huang, Zhang Zhang, Xiaotang Chen, Kaiqi Huang, and Shu Zhang. Towards rich feature discovery with class activation maps augmentation for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1389–1398, 2019.
- [106] Jing Xu, Rui Zhao, Feng Zhu, Huaming Wang, and Wanli Ouyang. Attention-aware compositional network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2119–2128, 2018.
- [107] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, Xin Jin, and Zhibo Chen. Relation-aware global attention for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3186–3195, 2020.
- [108] Hao Liu, Jiashi Feng, Meibin Qi, Jianguo Jiang, and Shuicheng Yan. End-to-end comparative attention networks for person re-identification. *IEEE Transactions on Image Processing*, 26(7):3492–3506, 2017.
- [109] Shang Gao, Jingya Wang, Huchuan Lu, and Zimo Liu. Pose-guided visible part matching for occluded person reid. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11744–11752, 2020.
- [110] Maoqing Tian, Shuai Yi, Hongsheng Li, Shihua Li, Xuesen Zhang, Jianping Shi, Junjie Yan, and Xiaogang Wang. Eliminating background-bias for robust person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5794–5803, 2018.
- [111] Thi Thanh Thuy Pham, Thi-Lan Le, Hai Vu, Trung Kien Dao, et al. Fully-automated person re-identification in multi-camera surveillance system with a robust kernel descriptor and effective shadow removal method. *Image and Vision Computing*, 59:44–62, 2017.
- [112] Xiang Bai, Mingkun Yang, Tengeng Huang, Zhiyong Dou, Rui Yu, and Yongchao Xu. Deep-person: Learning discriminative deep features for person re-identification. *Pattern Recognition*, 98:107036, 2020.
- [113] Shuai Li, Wenfeng Song, Zheng Fang, Jiaying Shi, Aimin Hao, Qingping Zhao, and Hong Qin. Long-short temporal-spatial clues excited network for robust person re-identification. *International Journal of Computer Vision*, 128(12):2936–2961, 2020.
- [114] Sanping Zhou, Jinjun Wang, Deyu Meng, Xiaomeng Xin, Yubing Li, Yihong Gong, and Nanning Zheng. Deep self-paced learning for person re-identification. *Pattern Recognition*, 76:739–751, 2018.
- [115] Sanping Zhou, Fei Wang, Zeyi Huang, and Jinjun Wang. Discriminative feature learning with consistent attention regularization for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8040–8049, 2019.
- [116] Sanping Zhou, Jinjun Wang, Deyu Meng, Yudong Liang, Yihong Gong, and Nanning Zheng. Discriminative feature learning with foreground attention for person re-identification. *IEEE Transactions on Image Processing*, 28(9):4671–4684, 2019.
- [117] Yiheng Liu, Wengang Zhou, Jianzhuang Liu, Guo-Jun Qi, Qi Tian, and Houqiang Li. An end-to-end foreground-aware network for person re-identification. *IEEE Transactions on Image Processing*, 30:2060–2071, 2021.
- [118] Jia Sun, Yanfeng Li, Houjin Chen, Bin Zhang, and Jinlei Zhu. Memf: Multi-level-attention embedding and multi-layer-feature fusion model for person re-identification. *Pattern Recognition*, 116:107937, 2021.
- [119] Xin Ning, Ke Gong, Weijun Li, Liping Zhang, Xiao Bai, and Shengwei Tian. Feature refinement and filter network for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [120] Yan Huang, Qiang Wu, Jingsong Xu, Yi Zhong, and Zhaoxiang Zhang. Unsupervised domain adaptation with background shift mitigating for person re-identification. *International Journal of Computer Vision*, 129(7):2244–2263, 2021.
- [121] Zhizheng Zhang, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Densely semantically aligned person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 667–676, 2019.
- [122] Lingxiao He, Jian Liang, Haiqing Li, and Zhenan Sun. Deep spatial feature reconstruction for partial person re-identification: Alignment-free approach. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7073–7082, 2018.
- [123] Yumin Suh, Jingdong Wang, Siyu Tang, Tao Mei, and Kyoung Mu Lee. Part-aligned bilinear representations for person re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 402–419, 2018.
- [124] Haiyu Zhao, Maoqing Tian, Shuyang Sun, Jing Shao, Junjie Yan, Shuai Yi, Xiaogang Wang, and Xiaoou 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*, pages 1077–1085, 2017.
- [125] Kuan Zhu, Haiyun Guo, Zhiwei Liu, Ming Tang, and Jinqiao Wang. Identity-guided human semantic parsing for person re-identification. *arXiv preprint arXiv:2007.13467*, 2020.
- [126] Yifan Sun, Qin Xu, Yali Li, Chi Zhang, Yikang Li, Shengjin Wang, and Jian Sun. Perceive where to focus: Learning visibility-aware part-level features for partial person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 393–402, 2019.
- [127] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [128] Mahdi M Kalayeh, Emrah Basaran, Muhittin Gökmen, Mustafa E Kamasak, and Mubarak Shah. Human semantic parsing for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1062–1071, 2018.
- [129] Hao Luo, Wei Jiang, Xuan Zhang, Xing Fan, Jingjing Qian, and Chi Zhang. Alignedreid++: Dynamically matching local information for person re-identification. *Pattern Recognition*, 94:53–61, 2019.
- [130] Rui Zhao, Wanli Ouyang, and Xiaogang Wang. Person re-identification by saliency learning. *IEEE transactions on pattern analysis and machine intelligence*, 39(2):356–370, 2016.

- [131] Sheng Li, Ming Shao, and Yun Fu. Person re-identification by cross-view multi-level dictionary learning. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2963–2977, 2017.
- [132] Changxing Ding, Kan Wang, Pengfei Wang, and Dacheng Tao. Multi-task learning with coarse priors for robust part-aware person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [133] Guanshuo Wang, Yufeng Yuan, Jiwei Li, Shiming Ge, and Xi Zhou. Receptive multi-granularity representation for person re-identification. *IEEE Transactions on Image Processing*, 29:6096–6109, 2020.
- [134] Ju Dai, Pingping Zhang, Dong Wang, Huchuan Lu, and Hongyu Wang. Video person re-identification by temporal residual learning. *IEEE Transactions on Image Processing*, 28(3):1366–1377, 2018.
- [135] Shuzhao Li, Huimin Yu, and Roland Hu. Attributes-aided part detection and refinement for person re-identification. *Pattern Recognition*, 97:107016, 2020.
- [136] Wei Shi, Hong Liu, and Mengyuan Liu. Image-to-video person re-identification using three-dimensional semantic appearance alignment and cross-modal interactive learning. *Pattern Recognition*, 122:108314, 2022.
- [137] Xiaolong Ma, Xiatian Zhu, Shaogang Gong, Xudong Xie, Jianming Hu, Kin-Man Lam, and Yisheng Zhong. Person re-identification by unsupervised video matching. *Pattern Recognition*, 65:197–210, 2017.
- [138] Slawomir Bak and Peter Carr. Deep deformable patch metric learning for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2690–2702, 2017.
- [139] Kan Liu, Bingpeng Ma, Wei Zhang, and Rui Huang. A spatio-temporal appearance representation for video-based pedestrian re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3810–3818, 2015.
- [140] Yang Shen, Weiyao Lin, Junchi Yan, Mingliang Xu, Jianxin Wu, and Jingdong Wang. Person re-identification with correspondence structure learning. In *Proceedings of the IEEE international conference on computer vision*, pages 3200–3208, 2015.
- [141] Zhong Zhang, Haijia Zhang, Shuang Liu, Yuan Xie, and Tariq S Durrani. Part-guided graph convolution networks for person re-identification. *Pattern Recognition*, 120:108155, 2021.
- [142] Tianyu He, Xin Jin, Xu Shen, Jianqiang Huang, Zhibo Chen, and Xian-Sheng Hua. Dense interaction learning for video-based person re-identification. *arXiv preprint arXiv:2103.09013*, 2021.
- [143] Pengfei Fang, Jieming Zhou, Soumava Kumar Roy, Lars Petersson, and Mehrtaash Harandi. Bilinear attention networks for person retrieval. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8030–8039, 2019.
- [144] Jianyuan Guo, Yuhui Yuan, Lang Huang, Chao Zhang, Jin-Ge Yao, and Kai Han. Beyond human parts: Dual part-aligned representations for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3642–3651, 2019.
- [145] Cheng Wang, Qian Zhang, Chang Huang, Wenyu Liu, and Xinggang Wang. Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 365–381, 2018.
- [146] Liming Zhao, Xi Li, Yueting Zhuang, and Jingdong Wang. Deeply-learned part-aligned representations for person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 3219–3228, 2017.
- [147] Fan Yang, Ke Yan, Shijian Lu, Huizhu Jia, Xiaodong Xie, and Wen Gao. Attention driven person re-identification. *Pattern Recognition*, 86:143–155, 2019.
- [148] Kan Wang, Changxing Ding, Stephen J Maybank, and Dacheng Tao. Cdpm: convolutional deformable part models for semantically aligned person re-identification. *IEEE Transactions on Image Processing*, 29:3416–3428, 2019.
- [149] Kan Wang, Pengfei Wang, Changxing Ding, and Dacheng Tao. Batch coherence-driven network for part-aware person re-identification. *IEEE Transactions on Image Processing*, 30:3405–3418, 2021.
- [150] Zhizhong Zhang, Yuan Xie, Ding Li, Wensheng Zhang, and Qi Tian. Learning to align via wasserstein for person re-identification. *IEEE Transactions on Image Processing*, 29:7104–7116, 2020.
- [151] Zhedong Zheng, Liang Zheng, and Yi Yang. Pedestrian alignment network for large-scale person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(10):3037–3045, 2018.
- [152] Xinqian Gu, Hong Chang, Bingpeng Ma, Hongkai Zhang, and Xilin Chen. Appearance-preserving 3d convolution for video-based person re-identification. In *European Conference on Computer Vision*, pages 228–243. Springer, 2020.
- [153] Houjing Huang, Wenjie Yang, Xiaotang Chen, Xin Zhao, Kaiqi Huang, Jinbin Lin, Guan Huang, and Dalong Du. Eanet: Enhancing alignment for cross-domain person re-identification. *arXiv preprint arXiv:1812.11369*, 2018.
- [154] Yanbei Chen, Xiatian Zhu, and Shaogang Gong. Person re-identification by deep learning multi-scale representations. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2590–2600, 2017.
- [155] Xuelin Qian, Yanwei Fu, Yu-Gang Jiang, Tao Xiang, and Xiangyang Xue. Multi-scale deep learning architectures for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5399–5408, 2017.
- [156] Feng Zheng, Cheng Deng, Xing Sun, Xinyang Jiang, Xiaowei Guo, Zongqiao Yu, Feiyue Huang, and Rongrong Ji. Pyramidal person re-identification via multi-loss dynamic training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8514–8522, 2019.
- [157] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3702–3712, 2019.
- [158] Jiawei Liu, Zheng-Jun Zha, Wei Wu, Kecheng Zheng, and Qibin Sun. Spatial-temporal correlation and topology learning for person re-identification in videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4370–4379, 2021.
- [159] Jianing Li, Shiliang Zhang, and Tiejun Huang. Multi-scale temporal cues learning for video person re-identification. *IEEE Transactions on Image Processing*, 29:4461–4473, 2020.
- [160] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Learning generalisable omni-scale representations for person re-identification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

- [161] Cairong Zhao, Xuekuan Wang, Wangmeng Zuo, Fumin Shen, Ling Shao, and Duoqian Miao. Similarity learning with joint transfer constraints for person re-identification. *Pattern Recognition*, 97:107014, 2020.
- [162] Ancong Wu, Wei-Shi Zheng, Xiaowei Guo, and Jian-Huang Lai. Distilled person re-identification: Towards a more scalable system. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1187–1196, 2019.
- [163] Yiluan Guo and Ngai-Man Cheung. Efficient and deep person re-identification using multi-level similarity. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2335–2344, 2018.
- [164] Cheng Yan, Guansong Pang, Lei Wang, Jile Jiao, Xuetao Feng, Chunhua Shen, and Jingjing Li. Bv-person: A large-scale dataset for bird-view person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10943–10952, 2021.
- [165] X Qian, Y Fu, T Xiang, YG Jiang, and X Xue. Leader-based multi-scale attention deep architecture for person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):371, 2020.
- [166] Wei Zhang, Xuanyu He, Xiaodong Yu, Weizhi Lu, Zhengjun Zha, and Qi Tian. A multi-scale spatial-temporal attention model for person re-identification in videos. *IEEE Transactions on Image Processing*, 29:3365–3373, 2019.
- [167] Xi Yang, Liangchen Liu, Nannan Wang, and Xinbo Gao. A two-stream dynamic pyramid representation model for video-based person re-identification. *IEEE Transactions on Image Processing*, 30:6266–6276, 2021.
- [168] Niki Martinel, Gian Luca Foresti, and Christian Micheloni. Deep pyramidal pooling with attention for person re-identification. *IEEE Transactions on Image Processing*, 29:7306–7316, 2020.
- [169] Guangyi Chen, Tianpei Gu, Jiwen Lu, Jin-An Bao, and Jie Zhou. Person re-identification via attention pyramid. *IEEE Transactions on Image Processing*, 30:7663–7676, 2021.
- [170] Yingji Zhong, Yaowei Wang, and Shiliang Zhang. Progressive feature enhancement for person re-identification. *IEEE Transactions on Image Processing*, 30:8384–8395, 2021.
- [171] Yewen Huang, Sicheng Lian, Haifeng Hu, Dihui Chen, and Tao Su. Multiscale omnibearing attention networks for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(5):1790–1803, 2020.
- [172] Hang Zhang, Huanhuan Cao, Xu Yang, Cheng Deng, and Dacheng Tao. Self-training with progressive representation enhancement for unsupervised cross-domain person re-identification. *IEEE Transactions on Image Processing*, 2021.
- [173] Xiaoxiao Sun and Liang Zheng. Dissecting person re-identification from the viewpoint of viewpoint. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 608–617, 2019.
- [174] Srikrishna Karanam, Yang Li, and Richard J Radke. Person re-identification with discriminatively trained viewpoint invariant dictionaries. In *Proceedings of the IEEE international conference on computer vision*, pages 4516–4524, 2015.
- [175] Shengcai Liao, Yang Hu, Xiangyu Zhu, and Stan Z Li. Person re-identification by local maximal occurrence representation and metric learning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2197–2206, 2015.
- [176] Angelo Porrello, Luca Bergamini, and Simone Calderara. Robust re-identification by multiple views knowledge distillation. In *European Conference on Computer Vision*, pages 93–110. Springer, 2020.
- [177] Ju Dai, Ying Zhang, Huchuan Lu, and Hongyu Wang. Cross-view semantic projection learning for person re-identification. *Pattern Recognition*, 75:63–76, 2018.
- [178] Ying-Cong Chen, Wei-Shi Zheng, Jian-Huang Lai, and Pong C Yuen. An asymmetric distance model for cross-view feature mapping in person reidentification. *IEEE transactions on circuits and systems for video technology*, 27(8):1661–1675, 2016.
- [179] Z Wu, Y Li, and RJ Radke. Viewpoint invariant human re-identification in camera networks using pose priors and subject-discriminative features. *IEEE transactions on pattern analysis and machine intelligence*, 37(5):1095, 2015.
- [180] Ying-Cong Chen, Xiatian Zhu, Wei-Shi Zheng, and Jian-Huang Lai. Person re-identification by camera correlation aware feature augmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(2):392–408, 2018.
- [181] Zhanxiang Feng, Jianhuang Lai, and Xiaohua Xie. Learning view-specific deep networks for person re-identification. *IEEE Transactions on Image Processing*, 27(7):3472–3483, 2018.
- [182] Shi-Zhe Chen, Chun-Chao Guo, and Jian-Huang Lai. Deep ranking for person re-identification via joint representation learning. *IEEE Transactions on Image Processing*, 25(5):2353–2367, 2016.
- [183] Jieru Jia, Qiuqi Ruan, Gaoyun An, and Yi Jin. Multiple metric learning with query adaptive weights and multi-task re-weighting for person re-identification. *Computer Vision and Image Understanding*, 160:87–99, 2017.
- [184] Lin Wu, Yang Wang, Zongyuan Ge, Qichang Hu, and Xue Li. Structured deep hashing with convolutional neural networks for fast person re-identification. *Computer Vision and Image Understanding*, 167:63–73, 2018.
- [185] Xiumei Chen, Xiangtao Zheng, and Xiaoqiang Lu. Bidirectional interaction network for person re-identification. *IEEE Transactions on Image Processing*, 30:1935–1948, 2021.
- [186] Alessandro Borgia, Yang Hua, Elyor Kodirov, and Neil M Robertson. Cross-view discriminative feature learning for person re-identification. *IEEE Transactions on Image Processing*, 27(11):5338–5349, 2018.
- [187] Jorge Garcia, Niki Martinel, Alfredo Gardel, Ignacio Bravo, Gian Luca Foresti, and Christian Micheloni. Discriminant context information analysis for post-ranking person re-identification. *IEEE Transactions on Image Processing*, 26(4):1650–1665, 2017.
- [188] Lin Wu, Yang Wang, Hongzhi Yin, Meng Wang, and Ling Shao. Few-shot deep adversarial learning for video-based person re-identification. *IEEE Transactions on Image Processing*, 29:1233–1245, 2020.
- [189] Houjing Huang, Wenjie Yang, Jinbin Lin, Guan Huang, Jiamiao Xu, Guoli Wang, Xiaotang Chen, and Kaiqi Huang. Improve person re-identification with part awareness learning. *IEEE Transactions on Image Processing*, 29:7468–7481, 2020.
- [190] Yutian Lin, Yu Wu, Chenggang Yan, Mingliang Xu, and Yi Yang. Unsupervised person re-identification via cross-camera similarity exploration. *IEEE Transactions on Image Processing*, 29:5481–5490, 2020.
- [191] Jingke Meng, Ancong Wu, and Wei-Shi Zheng. Deep asymmetric video-based person re-identification. *Pattern Recognition*, 93:430–441, 2019.

- [192] Cairong Zhao, Xuekuan Wang, Duoqian Miao, Hanli Wang, Weishi Zheng, Yong Xu, and David Zhang. Maximal granularity structure and generalized multi-view discriminant analysis for person re-identification. *Pattern Recognition*, 79:79–96, 2018.
- [193] Jieru Jia, Qiuqi Ruan, Yi Jin, Gaoyun An, and Shiming Ge. View-specific subspace learning and re-ranking for semi-supervised person re-identification. *Pattern Recognition*, 108:107568, 2020.
- [194] Hai-Miao Hu, Wen Fang, Bo Li, and Qi Tian. An adaptive multi-projection metric learning for person re-identification across non-overlapping cameras. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(9):2809–2821, 2019.
- [195] Lin Wu, Richang Hong, Yang Wang, and Meng Wang. Cross-entropy adversarial view adaptation for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(7):2081–2092, 2020.
- [196] Dapeng Tao, Yanan Guo, Baosheng Yu, Jianxin Pang, and Zhengtao Yu. Deep multi-view feature learning for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2657–2666, 2017.
- [197] Xueping Wang, Rameswar Panda, Min Liu, Yaonan Wang, and Amit K Roy-Chowdhury. Exploiting global camera network constraints for unsupervised video person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [198] Wei Zhang, Yimeng Li, Weizhi Lu, Xinshun Xu, Zhaowei Liu, and Xiangyang Ji. Learning intra-video difference for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 29(10):3028–3036, 2018.
- [199] Le An, Zhen Qin, Xiaojing Chen, and Songfan Yang. Multi-level common space learning for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(8):1777–1787, 2018.
- [200] Zhun Zhong, Liang Zheng, Zhiming Luo, Shaozi Li, and Yi Yang. Invariance matters: Exemplar memory for domain adaptive person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 598–607, 2019.
- [201] Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Cross-view asymmetric metric learning for unsupervised person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 994–1002, 2017.
- [202] Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Unsupervised person re-identification by deep asymmetric metric embedding. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):956–973, 2020.
- [203] Meng Zheng, Srikrishna Karanam, Ziyang Wu, and Richard J Radke. Re-identification with consistent attentive siamese networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5735–5744, 2019.
- [204] Lei Zhang, Fangyi Liu, and David Zhang. Adversarial view confusion feature learning for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(4):1490–1502, 2021.
- [205] Sicheng Lian, Weitao Jiang, and Haifeng Hu. Attention-aligned network for person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [206] Hao Chen, Yaohui Wang, Benoit Lagadec, Antitza Dantcheva, and Francois Bremond. Joint generative and contrastive learning for unsupervised person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2004–2013, 2021.
- [207] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8042–8051, 2018.
- [208] Xiang Li, Wei-Shi Zheng, Xiaojuan Wang, Tao Xiang, and Shaogang Gong. Multi-scale learning for low-resolution person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3765–3773, 2015.
- [209] Xiao-Yuan Jing, Xiaoke Zhu, Fei Wu, Xinge You, Qinglong Liu, Dong Yue, Ruimin Hu, and Baowen Xu. Super-resolution person re-identification with semi-coupled low-rank discriminant dictionary learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 695–704, 2015.
- [210] Ke Han, Yan Huang, Zerui Chen, Liang Wang, and Tieniu Tan. Prediction and recovery for adaptive low-resolution person re-identification. In *European Conference on Computer Vision*, pages 193–209. Springer, 2020.
- [211] Ke Han, Yan Huang, Chunfeng Song, Liang Wang, and Tieniu Tan. Adaptive super-resolution for person re-identification with low-resolution images. *Pattern Recognition*, page 107682, 2020.
- [212] XY Jing, X Zhu, F Wu, R Hu, X You, Y Wang, H Feng, and JY Yang. Super-resolution person re-identification with semi-coupled low-rank discriminant dictionary learning. *IEEE Transactions on Image Processing: a Publication of the IEEE Signal Processing Society*, 26(3):1363–1378, 2017.
- [213] Zhanxiang Feng, Jianhuang Lai, and Xiaohua Xie. Resolution-aware knowledge distillation for efficient inference. *IEEE Transactions on Image Processing*, 30:6985–6996, 2021.
- [214] Nikolaos Karianakis, Zicheng Liu, Yinpeng Chen, and Stefano Soatto. Reinforced temporal attention and split-rate transfer for depth-based person re-identification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 715–733, 2018.
- [215] Xiaokai Liu, Xiaorui Ma, Jie Wang, and Hongyu Wang. M3I: Multi-modality mining for metric learning in person re-identification. *Pattern Recognition*, 76:650–661, 2018.
- [216] Rahul Rama Varior, Gang Wang, Jiwen Lu, and Ting Liu. Learning invariant color features for person reidentification. *IEEE Transactions on Image Processing*, 25(7):3395–3410, 2016.
- [217] Zhiyi Cheng, Qi Dong, Shaogang Gong, and Xiatian Zhu. Inter-task association critic for cross-resolution person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2605–2615, 2020.
- [218] Yukun Huang, Zheng-Jun Zha, Xueyang Fu, Richang Hong, and Liang Li. Real-world person re-identification via degradation invariance learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14084–14094, 2020.
- [219] Zhedong Zheng, Xiaodong Yang, Zhiding Yu, Liang Zheng, Yi Yang, and Jan Kautz. Joint discriminative and generative learning for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2138–2147, 2019.
- [220] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3750–3759, 2019.

- [221] Xin Jin, Cuiling Lan, Wenjun Zeng, Zhibo Chen, and Li Zhang. Style normalization and restitution for generalizable person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3143–3152, 2020.
- [222] Weihua Chen, Xiaotang Chen, Jianguo Zhang, and Kaiqi Huang. Beyond triplet loss: a deep quadruplet network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 403–412, 2017.
- [223] Tong Xiao, Hongsheng Li, Wanli Ouyang, and Xiaogang Wang. Learning deep feature representations with domain guided dropout for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1249–1258, 2016.
- [224] Yunpeng Zhai, Qixiang Ye, Shijian Lu, Mengxi Jia, Rongrong Ji, and Yonghong Tian. Multiple expert brainstorming for domain adaptive person re-identification. *arXiv preprint arXiv:2007.01546*, 2020.
- [225] Zijie Zhuang, Longhui Wei, Lingxi Xie, Tianyu Zhang, Hengheng Zhang, Haozhe Wu, Haizhou Ai, and Qi Tian. Rethinking the distribution gap of person re-identification with camera-based batch normalization. In *European Conference on Computer Vision*, pages 140–157. Springer, 2020.
- [226] Chuanchen Luo, Chunfeng Song, and Zhaoxiang Zhang. Generalizing person re-identification by camera-aware invariance learning and cross-domain mixup. In *European Conference on Computer Vision*, volume 2, page 7. Springer, 2020.
- [227] Yang Zou, Xiaodong Yang, Zhiding Yu, BVK Kumar, and Jan Kautz. Joint disentangling and adaptation for cross-domain person re-identification. *arXiv preprint arXiv:2007.10315*, 2020.
- [228] Yan Bai, Jile Jiao, Wang Ce, Jun Liu, Yihang Lou, Xuetao Feng, and Ling-Yu Duan. Person30k: A dual-meta generalization network for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2123–2132, 2021.
- [229] Seokeon Choi, Taekyung Kim, Minki Jeong, Hyoungseob Park, and Changick Kim. Meta batch-instance normalization for generalizable person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3425–3435, 2021.
- [230] Yongxing Dai, Xiaotong Li, Jun Liu, Zekun Tong, and Ling-Yu Duan. Generalizable person re-identification with relevance-aware mixture of experts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16145–16154, 2021.
- [231] Anguo Zhang, Yueming Gao, Yuzhen Niu, Wenxi Liu, and Yongcheng Zhou. Coarse-to-fine person re-identification with auxiliary-domain classification and second-order information bottleneck. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 598–607, 2021.
- [232] Tianyu Zhang, Lingxi Xie, Longhui Wei, Zijie Zhuang, Yongfei Zhang, Bo Li, and Qi Tian. Unrealperson: An adaptive pipeline towards costless person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11506–11515, 2021.
- [233] Yi Zheng, Shixiang Tang, Guolong Teng, Yixiao Ge, Kaijian Liu, Jing Qin, Donglian Qi, and Dapeng Chen. Online pseudo label generation by hierarchical cluster dynamics for adaptive person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8371–8381, 2021.
- [234] Amena Khatun, Simon Denman, Sridha Sridharan, and Clinton Fookes. Joint identification–verification for person re-identification: A four stream deep learning approach with improved quartet loss function. *Computer Vision and Image Understanding*, 197:102989, 2020.
- [235] Amena Khatun, Simon Denman, Sridha Sridharan, and Clinton Fookes. A deep four-stream siamese convolutional neural network with joint verification and identification loss for person re-detection. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1292–1301. IEEE, 2018.
- [236] Hao Feng, Minghao Chen, Jinming Hu, Dong Shen, Haifeng Liu, and Deng Cai. Complementary pseudo labels for unsupervised domain adaptation on person re-identification. *IEEE Transactions on Image Processing*, 30:2898–2907, 2021.
- [237] Wenfeng Song, Shuai Li, Tao Chang, Aimin Hao, Qiping Zhao, and Hong Qin. Context-interactive cnn for person re-identification. *IEEE Transactions on Image Processing*, 29:2860–2874, 2019.
- [238] Yongxing Dai, Jun Liu, Yan Bai, Zekun Tong, and Ling-Yu Duan. Dual-refinement: Joint label and feature refinement for unsupervised domain adaptive person re-identification. *IEEE Transactions on Image Processing*, 30:7815–7829, 2021.
- [239] Shan Lin, Chang-Tsun Li, and Alex C Kot. Multi-domain adversarial feature generalization for person re-identification. *IEEE Transactions on Image Processing*, 30:1596–1607, 2020.
- [240] Hengheng Zhang, Ying Li, Zijie Zhuang, Lingxi Xie, and Qi Tian. 3d-gat: 3d-guided adversarial transform network for person re-identification in unseen domains. *Pattern Recognition*, 112:107799, 2021.
- [241] Zhicheng Zhao, Binlin Zhao, and Fei Su. Person re-identification via integrating patch-based metric learning and local salience learning. *Pattern Recognition*, 75:90–98, 2018.
- [242] Huafeng Li, Zhenyu Kuang, Zhengtao Yu, and Jiebo Luo. Structure alignment of attributes and visual features for cross-dataset person re-identification. *Pattern Recognition*, 106:107414, 2020.
- [243] Qize Yang, Hong-Xing Yu, Ancong Wu, and Wei-Shi Zheng. Patch-based discriminative feature learning for unsupervised person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3633–3642, 2019.
- [244] Hantao Yao, Shiliang Zhang, Richang Hong, Yongdong Zhang, Changsheng Xu, and Qi Tian. Deep representation learning with part loss for person re-identification. *IEEE Transactions on Image Processing*, 28(6):2860–2871, 2019.
- [245] Chi Su, Shiliang Zhang, Junliang Xing, Wen Gao, and Qi Tian. Deep attributes driven multi-camera person re-identification. In *European conference on computer vision*, pages 475–491. Springer, 2016.
- [246] Zhiyuan Shi, Timothy M Hospedales, and Tao Xiang. Transferring a semantic representation for person re-identification and search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4184–4193, 2015.
- [247] Peixi Peng, Tao Xiang, Yaowei Wang, Massimiliano Pontil, Shaogang Gong, Tiejun Huang, and Yonghong Tian. Unsupervised cross-dataset transfer learning for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1306–1315, 2016.
- [248] Xin Jin, Cuiling Lan, Wenjun Zeng, and Zhibo Chen. Global distance-distributions separation for unsupervised person re-identification. In *European Conference on Computer Vision*, pages 735–751. Springer, 2020.

- [249] Guangyi Chen, Yuhao Lu, Jiwen Lu, and Jie Zhou. Deep credible metric learning for unsupervised domain adaptation person re-identification. In *Proc. Eur. Conf. Comput. Vis.*, pages 643–659. Springer, 2020.
- [250] Djebri Mekhazni, Amran Bhuiyan, George Ekladios, and Eric Granger. Unsupervised domain adaptation in the dissimilarity space for person re-identification. In *European Conference on Computer Vision*, pages 159–174. Springer, 2020.
- [251] Zechen Bai, Zhigang Wang, Jian Wang, Di Hu, and Errui Ding. Unsupervised multi-source domain adaptation for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12914–12923, 2021.
- [252] Dengpan Fu, Dongdong Chen, Jianmin Bao, Hao Yang, Lu Yuan, Lei Zhang, Houqiang Li, and Dong Chen. Unsupervised pre-training for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14750–14759, 2021.
- [253] Shiyu Xuan and Shiliang Zhang. Intra-inter camera similarity for unsupervised person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11926–11935, 2021.
- [254] Fengxiang Yang, Zhun Zhong, Zhiming Luo, Yuanzheng Cai, Yaojin Lin, Shaozi Li, and Nicu Sebe. Joint noise-tolerant learning and meta camera shift adaptation for unsupervised person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4855–4864, 2021.
- [255] Xiao Zhang, Yixiao Ge, Yu Qiao, and Hongsheng Li. Refining pseudo labels with clustering consensus over generations for unsupervised object re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3436–3445, 2021.
- [256] Yuyang Zhao, Zhun Zhong, Fengxiang Yang, Zhiming Luo, Yaojin Lin, Shaozi Li, and Nicu Sebe. Learning to generalize unseen domains via memory-based multi-source meta-learning for person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6277–6286, 2021.
- [257] Kecheng Zheng, Wu Liu, Lingxiao He, Tao Mei, Jiebo Luo, and Zheng-Jun Zha. Group-aware label transfer for domain adaptive person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5310–5319, 2021.
- [258] Takashi Isobe, Dong Li, Lu Tian, Weihua Chen, Yi Shan, and Shengjin Wang. Towards discriminative representation learning for unsupervised person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8526–8536, 2021.
- [259] Yu Zhao, Qiaoyuan Shu, Keren Fu, Pengcheng Wei, and Jian Zhan. Joint patch and instance discrimination learning for unsupervised person re-identification. *Image and Vision Computing*, 103:104000, 2020.
- [260] Yan Bai, Ce Wang, Yihang Lou, Jun Liu, and Ling-Yu Duan. Hierarchical connectivity-centered clustering for unsupervised domain adaptation on person re-identification. *IEEE Transactions on Image Processing*, 30:6715–6729, 2021.
- [261] Meng Liu, Leigang Qu, Liqiang Nie, Maofu Liu, Lingyu Duan, and Baoquan Chen. Iterative local-global collaboration learning towards one-shot video person re-identification. *IEEE Transactions on Image Processing*, 29:9360–9372, 2020.
- [262] Kongzhu Jiang, Tianzhu Zhang, Yongdong Zhang, Feng Wu, and Yong Rui. Self-supervised agent learning for unsupervised cross-domain person re-identification. *IEEE Transactions on Image Processing*, 29:8549–8560, 2020.
- [263] Jia Sun, Yanfeng Li, Houjin Chen, Yahui Peng, and Jinlei Zhu. Unsupervised cross domain person re-identification by multi-loss optimization learning. *IEEE Transactions on Image Processing*, 30:2935–2946, 2021.
- [264] Minxian Li, Xiatian Zhu, and Shaogang Gong. Unsupervised tracklet person re-identification. *IEEE transactions on pattern analysis and machine intelligence*, 42(7):1770–1782, 2020.
- [265] Huafeng Li, Shuanglin Yan, Zhengtao Yu, and Dapeng Tao. Attribute-identity embedding and self-supervised learning for scalable person re-identification. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(10):3472–3485, 2019.
- [266] Jingke Meng, Sheng Wu, and Wei-Shi Zheng. Weakly supervised person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 760–769, 2019.
- [267] Yunpeng Zhai, Shijian Lu, Qixiang Ye, Xuebo Shan, Jie Chen, Rongrong Ji, and Yonghong Tian. Ad-cluster: Augmented discriminative clustering for domain adaptive person re-identification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9021–9030, 2020.
- [268] Shengcai Liao and Ling Shao. Interpretable and generalizable person re-identification with query-adaptive convolution and temporal lifting. *arXiv preprint arXiv:1904.10424*, 2019.
- [269] Sanping Zhou, Jinjun Wang, Jiayun Wang, Yihong Gong, and Nanning Zheng. Point to set similarity based deep feature learning for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3741–3750, 2017.
- [270] Yingzhi Tang, Xi Yang, Nannan Wang, Bin Song, and Xinbo Gao. Cgan-tm: A novel domain-to-domain transferring method for person re-identification. *IEEE Transactions on Image Processing*, 29:5641–5651, 2020.
- [271] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 994–1003, 2018.
- [272] Nan Pu, Wei Chen, Yu Liu, Erwin M Bakker, and Michael S Lew. Lifelong person re-identification via adaptive knowledge accumulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7901–7910, 2021.
- [273] Liangchen Song, Cheng Wang, Lefei Zhang, Bo Du, Qian Zhang, Chang Huang, and Xinggang Wang. Unsupervised domain adaptive re-identification: Theory and practice. *Pattern Recognition*, 102:107173, 2020.
- [274] Binghui Chen, Weihong Deng, and Jiani Hu. Mixed high-order attention network for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 371–381, 2019.
- [275] Chiat-Pin Tay, Sharmili Roy, and Kim-Hui Yap. Aanet: Attribute attention network for person re-identifications. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7134–7143, 2019.
- [276] Jianlou Si, Honggang Zhang, Chun-Guang Li, Jason Kuen, Xiangfei Kong, Alex C Kot, and Gang Wang. Dual attention matching network for context-aware feature sequence based person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5363–5372, 2018.

- [277] Jieming Zhou, Soumava Kumar Roy, Pengfei Fang, Mehrtash Harandi, and Lars Petersson. Cross-correlated attention networks for person re-identification. *Image and Vision Computing*, 100:103931, 2020.
- [278] Wei Li, Xiatian Zhu, and Shaogang Gong. Scalable person re-identification by harmonious attention. *International Journal of Computer Vision*, 128(6):1635–1653, 2020.
- [279] Dengpan Fu, Bo Xin, Jingdong Wang, Dongdong Chen, Jianmin Bao, Gang Hua, and Houqiang Li. Improving person re-identification with iterative impression aggregation. *IEEE Transactions on Image Processing*, 29:9559–9571, 2020.
- [280] Shuangjie Xu, Yu Cheng, Kang Gu, Yang Yang, Shiyu Chang, and Pan Zhou. Jointly attentive spatial-temporal pooling networks for video-based person re-identification. In *Proceedings of the IEEE international conference on computer vision*, pages 4733–4742, 2017.