

A Systematic Evaluation and Benchmark for Person Re-Identification: Features, Metrics, and Datasets

Srikrishna Karanam*, *Student Member, IEEE*, Mengran Gou*, *Student Member, IEEE*,
Ziyan Wu, *Member, IEEE*, Angels Rates-Borras, Octavia Camps, *Member, IEEE*,
and Richard J. Radke, *Senior Member, IEEE*

Abstract—Person re-identification (re-id) is a critical problem in video analytics applications such as security and surveillance. The public release of several datasets and code for vision algorithms has facilitated rapid progress in this area over the last few years. However, directly comparing re-id algorithms reported in the literature has become difficult since a wide variety of features, experimental protocols, and evaluation metrics are employed. In order to address this need, we present an extensive review and performance evaluation of single- and multi-shot re-id algorithms. The experimental protocol incorporates the most recent advances in both feature extraction and metric learning. To ensure a fair comparison, all of the approaches were implemented using a unified code library that includes 8 feature extraction algorithms and 19 metric learning and ranking techniques. All approaches were evaluated using a new large-scale dataset that closely mimics a real-world problem setting, in addition to 13 other publicly available datasets: VIPeR, GRID, CAVIAR, 3DPeS, PRID, V47, WARD, SAIPT-SoftBio, CUHK03, RAiD, iLIDSVID, HDA+ and Market1501. The evaluation codebase and results will be made publicly available for community use.



1 INTRODUCTION

PERSON re-identification, or re-id, is a critical task in most surveillance and security applications [1], [2], [3] and has increasingly attracted attention from the computer vision community [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21], [22], [23]. The fundamental re-id problem is to compare a person of interest as seen in a “probe” camera view to a “gallery” of candidates captured from a camera that does not overlap with the probe one. If a true match to the probe exists in the gallery, it should have a high matching score, or rank, compared to incorrect candidates.

Since the body of research in re-id is now quite large, we can begin to draw conclusions about the best combinations of algorithmic subcomponents. In this paper, we present a careful, fair, and systematic evaluation of feature extraction, metric learning, and multi-shot ranking algorithms proposed for re-id on a wide variety of benchmark datasets. Our general evaluation framework is to consider all possible combinations of feature extraction and metric learning

algorithms for single-shot datasets and all possible combinations of feature extraction, metric learning, and multi-shot ranking algorithms for multi-shot datasets. In particular, we evaluate 208 such algorithm combinations on 7 single-shot re-id datasets and 448 such algorithm combinations on 7 multi-shot re-id datasets, making the proposed study the **largest and most systematic** re-id benchmark to date. As part of the evaluation, we built a **public code library** with an easy-to-use input/output code structure and uniform algorithm parameters that includes 8 contemporary feature extraction and 19 metric learning and ranking algorithms. Both the code library and the complete benchmark results will be made publicly available for community use.

Existing re-id algorithms are typically evaluated on academic re-id datasets [4], [24], [25], [26], [27], [28], [29], [30] that are specifically hand-curated to only have sets of bounding boxes for the probes and the corresponding matching candidates. On the other hand, real-world end-to-end surveillance systems include automatic detection and tracking modules, depicted in Figure 1, that generate candidates on-the-fly, resulting in gallery sets that are dynamic in nature. Furthermore, errors in these modules may result in bounding boxes that may not accurately represent a human [3]. While these issues are critical in practical re-id applications, they are not well-represented in the currently available datasets. To this end, our evaluation also includes a **new, large-scale dataset** constructed from images captured in a challenging surveillance camera network from an airport. All the images in this dataset were generated by running a prototype end-to-end real-time re-id system using automatic person detection and tracking algorithms instead of hand-curated bounding boxes.

- S. Karanam and R.J. Radke are with the Department of Electrical, Computer, and Systems Engineering, Rensselaer Polytechnic Institute, Troy, NY, 12180 (e-mail: rjradke@ecse.rpi.edu).
- Z. Wu is with Siemens Corporation, Corporate Technology, Princeton, NJ 08540 (e-mail: ziyuan.wu@siemens.com).
- M. Gou, A. Rates-Borras, and O. Camps are with the Department of Electrical and Computer Engineering, Northeastern University, Boston, MA, 02115 (e-mail: camps@coe.neu.edu).
- *The first two authors contributed equally to this work. Corresponding author: S. Karanam.

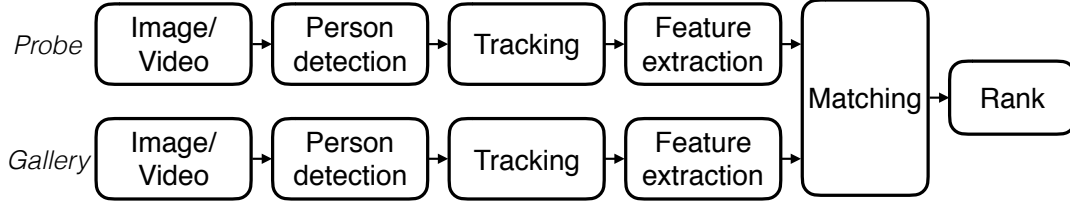


Fig. 1. A typical end-to-end re-id system pipeline.

2 EVALUATED TECHNIQUES

In this section, we summarize the feature extraction, metric learning, and multi-shot ranking techniques that are evaluated as part of the proposed re-id benchmark, which include algorithms published through CVPR 2016. We anticipate that the benchmark will be updated (along the lines of the Middlebury benchmarks [31], [32]) as new techniques are implemented into our evaluation framework.

2.1 Feature extraction

We consider 8 feature extraction schemes that are commonly used in the re-id literature, summarized in Table 1(a). In ELF [4], color histograms in the RGB, YCbCr, and HS color spaces, and texture histograms of responses of rotationally invariant Schmid [33] and Gabor [34] filters are computed. In LDFV [35], local pixel descriptors comprising pixel spatial location, intensity, and gradient information are encoded into the Fisher vector [36] representation. In gBiCov [37], multi-scale biologically-inspired features [38] are encoded using covariance descriptors [39]. In AlexNet-Finetune, we start with a convolutional neural network (CNN) based on the AlexNet architecture [40] pre-trained on the ImageNet dataset, and finetune it using the datasets we consider in this evaluation. Specifically, during finetuning, we do not modify the weights of the convolutional layers and train the last two fully-connected layers from scratch. More implementation details are presented in Section 3.2. In DenseColorSIFT [9], each image is densely divided into patches, and color histograms and SIFT features are extracted from each patch. In HistLBP [13], color histograms in the RGB, YCbCr, and HS color spaces and texture histograms from local binary patterns (LBP) [41] features are computed. In LOMO [18], HSV color histograms and scale-invariant LBP [42] features are extracted from the image processed by a multi-scale Retinex algorithm [43], and maximally-pooled along the same horizontal strip. In GOG [44], an image is divided into horizontal strips and local patches in each strip are modeled using a Gaussian distribution. Each strip is then regarded as a set of such Gaussian distributions, which is then summarized using a single Gaussian distribution.

2.2 Metric learning

While using any of the features described in the previous section in combination with the Euclidean distance (l_2) can be used to rank gallery candidates, this would be an unsupervised and suboptimal approach. Incorporating supervision using training data leads to superior performance, which is the goal of metric learning, i.e., learning a new feature space such that feature vectors of the same

person are close whereas those of different people are relatively far. We consider 16 metric learning methods that are typically used by the re-id community, summarized in Table 1(b). Fisher discriminant analysis (FDA) [45], local Fisher discriminant analysis (LFDA) [11], marginal Fisher analysis (MFA) [47], and cross-view quadratic discriminant analysis (XQDA) [18] all formulate a Fisher-type optimization problem that seeks to minimize the within-class data scatter while maximizing between-class data scatter. In practice, scatter matrices are regularized by a small fraction of their trace to deal with matrix singularities. Information-theoretic metric learning (ITML) [46], large-margin nearest neighbor (LMNN) [49], relative distance comparison (PRDC) [6], keep-it-simple-and-straightforward metric (KISSME) [7], and pairwise constrained component analysis (PCCA) [8] all learn Mahalanobis-type distance functions using variants of the basic pairwise constraints principle. kPCCA [8], kLFDA [13], and kMFA [13] kernelize PCCA, LFDA, and MFA, respectively. For these kernel-based methods, we consider the standard linear, exponential (exp), chi2 (χ^2), and chi2-rbf (R_{χ^2}) kernels. In RankSVM [5], a weight vector that weights the different features appropriately is learned using a soft-margin SVM formulation. In SVMML [48], locally adaptive decision functions are learned in a large-margin SVM framework.

2.3 Multi-shot ranking

While most re-id algorithms are single-shot, i.e. features are extracted from a single probe image of the person of interest, the multi-shot scenario, in which features are extracted from a series of images of the person of interest, is arguably more relevant to video analysis problems. The simplest way to handle multi-shot data is to compute the average feature vector for each person, effectively resulting in a single-shot problem. However, we also evaluated several algorithms that inherently address multi-shot data, treating it as an image set and constructing affine hulls to compute the distance between a gallery and a probe person. Specifically, we considered the AHISD [50] and RNP [51] algorithms. While these methods were proposed in the context of face recognition, the basic notion of image set matching applies to re-id as well. We also evaluated a multi-shot method based on sparse ranking, in which re-id is posed as a sparse recovery problem. Specifically, we considered SRID [52], where a block sparse recovery problem is solved to retrieve the identity of a probe person.

2.4 Techniques not (yet) considered

As noted in Section 1, the framework we adopt involves evaluating all possible combinations of candidate feature extraction, metric learning, and multi-shot ranking algorithms.

Feature	Year	Metric	Year	Metric	Year
ELF [4]	ECCV 08	l_2		PCCA [8]	CVPR 12
LDFV [35]	ECCVW 12	FDA [45]	AE 1936	kPCCA [8]	CVPR 12
gBiCov [37]	BMVC 12	ITML [46]	ICML 07	LFDA [11]	CVPR 13
AlexNet-Finetune [40]	NIPS 12	MFA [47]	PAMI 07	SVMML [48]	CVPR 13
DenseColorSIFT [9]	CVPR 13	LMNN [49]	JMLR 08	kMFA [13]	ECCV 14
HistLBP [13]	ECCV 14	RankSVM [5]	BMVC 10	rPCCA [13]	ECCV 14
LOMO [18]	CVPR 15	PRDC [6]	CVPR 11	kLFDA [13]	ECCV 14
GOG [44]	CVPR 16	KISSME [7]	CVPR 12	XQDA [18]	CVPR 15

(a)
(b)

TABLE 1
Evaluated feature extraction and metric learning methods.

Methods that do not fall into this evaluation framework include post-rank learning methods [53], [54], unsupervised learning [9], [24], [55], attribute learning [56], [57], [58], ensemble methods [59], [60], [61] and mid-level representation learning [14]. A more comprehensive survey of these and other related methods can be found in the book by Gong *et al.* [2] and papers by Zheng [62], Satta [63], Vezzani [64], and Bedagkar-Gala and Shah [65]. While these methods are currently not part of our evaluation, we plan to expand our study and include them in a future release.

3 DATASETS

In this section, we briefly summarize the various publicly available datasets that are used in our benchmark evaluation. Table 2 provides a statistical summary of each dataset. Based on difficult examples, we also annotate each dataset with challenging attributes from the following list: viewpoint variations (VV), illumination variations (IV), detection errors (DE), occlusions (OCC), background clutter (BC), and low-resolution images (RES). We also indicate the number of bounding boxes (BBox) and cameras (cam) in each dataset, and the means by which the bounding boxes were obtained: using hand-labeling (hand), aggregated channel features [66] (ACF), or the deformable parts model detector [67] (DPM). A few sample difficult examples for both single- and multi-shot datasets are shown in Figures 2 and 3. More examples are provided as part of the supplementary material.

VIPeR [4] consists of 632 people from two disjoint views. Each person has only one image per view. As can be noted from Figure 2(c), VIPeR suffers from viewpoint and illumination variations. GRID [68] has 250 image pairs collected from 8 non-overlapping cameras. To mimic a realistic scenario, 775 non-paired people are included in the gallery set, which makes this dataset extremely challenging. GRID suffers from viewpoint variations, background clutter, occlusions and low-resolution images (Figure 2(d)). CAVIAR [24] is constructed from two cameras in a shopping mall. Of the 72 people, we only use 50 people who have images from both cameras. CAVIAR suffers from viewpoint variations

and low-resolution images. In the case of 3DPeS [26], the re-id community uses a set of selected snapshots instead of the original video, which includes 192 people and 1,011 images. 3DPeS suffers from viewpoint and illumination variations (Figure 2(a)). PRID [25] is constructed from two outdoor cameras, with 385 tracking sequences from one camera and 749 tracking sequences from the other camera. Among them, only 200 people appear in both views. To be consistent with previous work [28], we use the same subset of the data with 178 people. PRID suffers from viewpoint and illumination variations (Figure 3(b)). V47 [69] contains 47 people walking through two indoor cameras. WARD [27] collects 4,786 images of 70 people in 3 disjoint cameras. SAIVT-Softbio [70] consists of 152 people as seen from a surveillance camera network with 8 cameras. To be consistent with existing work [70], we use only two camera pairs: camera 3 and camera 8 (which we name SAIVT-38) and camera 5 and camera 8 (which we name SAIVT-58). Both these datasets suffer from viewpoint and illumination variations, and background clutter (Figure 3(c)). CUHK03 [71] has 1360 people and 13,164 images from 5 disjoint camera pairs. Both manually labeled bounding boxes and automatically detected bounding boxes using the DPM detector [67] are provided. We only use the detected bounding boxes in our experiments. CUHK03 suffers from viewpoint variations, detection errors, and occlusions. RAiD [29] includes 43 people as seen from two indoor and two outdoor cameras and suffers from viewpoint and illumination variations. iLIDSVID [28] includes 600 tracking sequences for 300 people from 2 non-overlapping cameras in an airport and suffers from viewpoint and illumination variations, background clutter, and occlusions (Figure 3(a)). HDA+ [72] was proposed to be a testbed for an automatic re-id system. Fully labeled frames for 30-minute long videos from 13 disjoint cameras are provided. Since we only focus on the re-id problem, we use pre-detected bounding boxes generated using the ACF [66] detector. Market1501 [73] has 1,501 people with 32,643 images and 2,793 false alarms from a person detector [67]. Besides these, an additional 500,000 false alarms and non-paired people are also provided to emphasize practical problems in re-id. Market1501 suffers from viewpoint variations, detection errors and low-resolution



Fig. 2. Difficult examples for single-shot datasets: (a) 3DPeS, (b) Airport, (c) VIPeR, and (d) GRID. The first and second rows show the probe and gallery images respectively and the number in third row is the corresponding rank from $kLFDA_l$ with GOG.

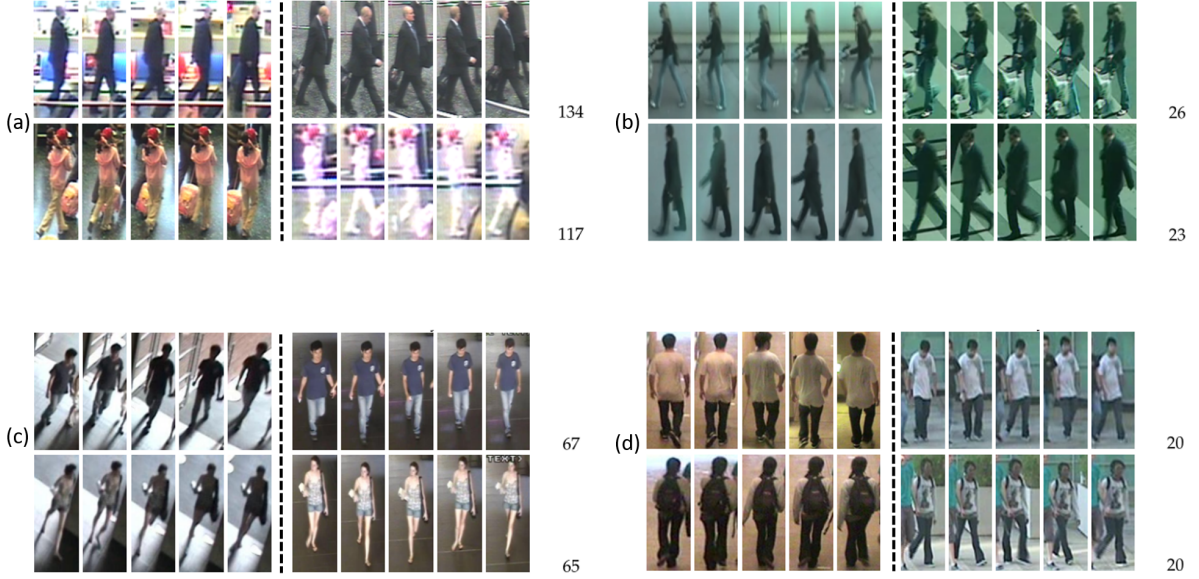


Fig. 3. Difficult examples for multi-shot datasets: (a) iLIDSVID, (b) PRID, (c) SAIVT-58, and (d) WARD-13. In each image, images on the left of the dividing line show the set of images for a probe person, and the images on the right show set of images of the same person in the gallery view. The number towards the end is the corresponding rank from KISSME and SRID with GOG.

(Figure 2(e)). Airport is the new dataset we introduce in the next section.

3.1 A new, real-world, large-scale dataset

In this section, we provide details about a new re-id dataset we designed for this benchmark. The dataset was created using videos from six cameras of an indoor surveillance network in a mid-sized airport; this testbed is described further in [3]. The cameras cover various parts of a central security checkpoint area and three concourses. Each camera has 768×432 pixels and captures video at 30 frames per second. 12-hour long videos from 8 AM to 8 PM were collected from each of these cameras. Under the assumption that each target person takes a limited amount of time to travel through the network, each of these long videos was randomly split into 40 five minute long video clips. Each video clip was then run through a prototype end-to-end re-id system comprised of automatic person detection and tracking algorithms. Specifically, we employed the ACF framework of Dollar *et al.* [66] to detect people and a combi-

nation of FAST corner features [74] and the KLT tracker [75] to track people and associate any broken “tracklets”.

Since all the bounding boxes were generated automatically without any manual annotation, this dataset accurately mimics a real-world re-id problem setting. A typical fully automatic re-id system should be able to automatically detect, track, and match people seen in the gallery camera, and the proposed dataset exactly reflects this setup. In total, from all the short video clips, tracks corresponding to 9,651 unique people were extracted. The number of bounding box images in the dataset is 39,902, giving an average of 3.13 images per person. The sizes of detected bounding boxes range from 130×54 to 403×166 . 1,382 of the 9,651 people are paired in at least two cameras. A number of unpaired people are also included in the dataset to simulate how a real-world re-id system would work: given a person of interest in the probe camera, a real system would automatically detect and track all the people seen in the gallery camera. Therefore, having a dataset with a large number of unpaired people greatly facilitates algorithmic re-id research by closely simulating a real-world environment. A sample of

Dataset	# people	# BBox	# distractors	# cam	label	Attributes
VIPeR	632	1,264	0	2	hand	VV,IV
GRID	250	500	775	2	hand	VV,BC,OCC,RES
CAVIAR	72	1,220	0	2	hand	VV,RES
3DPeS	192	1,011	0	8	hand	VV,IV
PRID	178	24,541	0	2	hand	VV,IV
V47	47	752	0	2	hand	-
WARD	70	4,786	0	3	hand	IV
SAIVT-Softbio	152	64,472	0	8	hand	VV,IV,BC
CUHK03	1,360	13,164	0	2	DPM/hand	VV,DE,OCC
RAiD	43	6,920	0	4	hand	VV,IV
iLIDSVID	300	42,495	0	2	hand	VV,IV,BC,OCC
HDA+	74	2,976	0	12	ACF/hand	VV,IV,DE
Market1501	1,501	32,643	2,793+500K	6	DPM	VV,DE,RES
Airport	1,382	8,664	31,238	6	ACF	VV,IV,DE,BC,OCC

TABLE 2
The characteristics of the 14 datasets of the re-id benchmark.

the images available in the dataset is shown in Figure 4. As can be seen from the figure, these are the kind of images one would expect from a fully automated system with detection and tracking modules working in a real-world surveillance environment. As shown in Figure 2(b), the Airport dataset suffers from all challenging attributes except low resolution. That is because relatively small detections are rejected by the person detector.

3.2 Evaluation protocol

3.2.1 Datasets, and training and testing splits.

Based on the number of images for each probe person, we categorize the datasets into either the single-shot or multi-shot setting. We employ the single-shot evaluation protocol for VIPeR, GRID, 3DPeS, CUHK03, HDA+, Market1501 and Airport. For the other 7 datasets, we employ the multi-shot evaluation protocol. In the Airport dataset, we fix one of the 6 cameras as the probe view and randomly pick paired people from 20 of the 40 short clips as the training set, with the rest forming the testing set. In the case of GRID, CUHK03, HDA+ and Market1501, we use the partition files provided by the respective authors. In RAiD, we fix camera 1 as the probe view, resulting in three sub-datasets, RAiD-12, RAiD-13, and RAiD-14, corresponding to the 3 possible gallery views. RAiD-12 has 43 people, of which we use 21 people to construct the training set and the rest to construct the testing set. The other two sub-datasets have 42 people each, which we split into equal-sized training and testing sets. In WARD, we fix camera 1 as the probe view, resulting in two sub-datasets, WARD-12 and WARD-13, corresponding to the 2 possible gallery views. Both these sub-datasets have 70 people each. We split VIPeR, GRID, CAVIAR, 3DPeS, PRID, WARD-12, WARD-13 and iLIDSVID into equal-sized training and testing sets. SAIVT-38 has 99 people, of which we use 31 people for training and the rest for testing. SAIVT-58 has 103 people, of which we use 33 people for training and the rest for testing. Finally, for each

dataset, we use 10 different randomly generated training and testing sets and report the overall average results.

3.2.2 Evaluation framework.

In the single-shot evaluation scheme, for each dataset, we consider two aspects: type of feature and type of metric learning algorithm. We evaluate all possible combinations of the 8 different features and 16 different metric learning algorithms listed in Table 1. Since we also evaluate four different kernels for the kernelized algorithms, the total number of algorithm combinations is 208.¹ In the multi-shot evaluation scheme, we consider three aspects: type of feature, type of metric learning algorithm, and type of ranking algorithm. Additionally, we consider two evaluation sub-schemes: using the average feature vector as the data representative (called AVER), and clustering the multiple feature vectors for each person and considering the resulting cluster centers as the representative feature vectors for each person (called CLUST). AVER effectively converts each dataset into an equivalent single-shot version, giving the same 208 algorithm combinations as above. However, in the case of CLUST, we do not consider kernelized metric learning algorithms and other non-typical algorithms such as RankSVM and SVMML because only AVER can be employed to rank gallery candidates. Consequently, we use the remaining 9 metric learning algorithms and the baseline l_2 method, in which we use the features in the original space without any projection. These 10 algorithms are used in combination with the 8 different features and 3 different ranking algorithms, resulting in 240 algorithm combinations for CLUST. Therefore, in total, we evaluate 448 different algorithm combinations for each multi-shot dataset.

1. Since LDFV and GOG features are not non-negative, we evaluate only linear and exp kernels in this case.

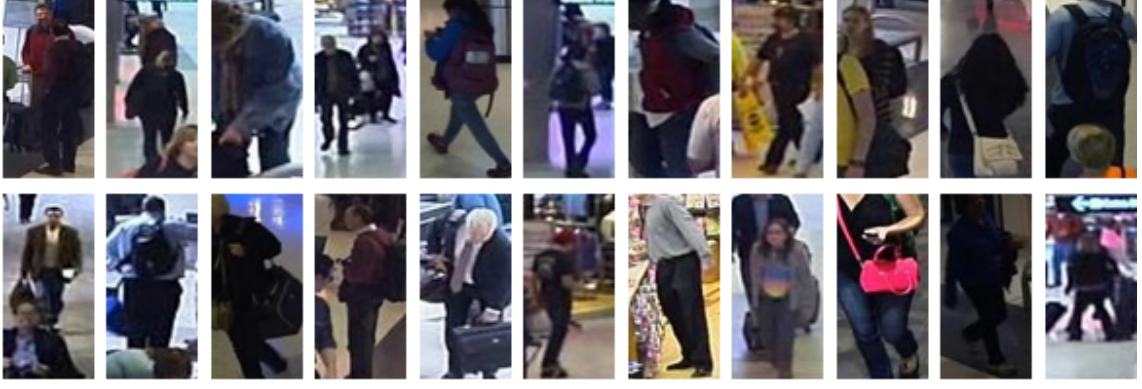


Fig. 4. Samples of images from the proposed Airport dataset.

3.2.3 Implementation and parameter details.

We normalize all images of a particular dataset to the same size, which is set to 128×48 for VIPeR, GRID, CAVIAR and 3DPeS and 128×64 for all other datasets. To compute features, we divide each image into 6 horizontal rectangular strips. In the case of LOMO, since the patches are fixed to be square-shaped, we obtain 12 patches for a 128×48 image and 18 patches for a 128×64 image.

In the case of AlexNet-Finetune, we resize each image to 227×227 pixels following [40]. We start training with a model pre-trained on the ImageNet dataset and train the fully connected layers fc7 and fc8 from scratch. The number of output units in the fc7 layer is set to 4096. Since we consider each person to be a different class, we set the number of output units in the fc8 layer to the number of unique people in our training set. Depending on the training split, this number varies from 2560 to 2580. Once the model is trained, we use the output of the fc7 layer as the image descriptor, giving a 4096-dimensional feature vector.

In metric learning, we set the projected feature space dimension to 40. We set the ratio of the number of negative to positive pairwise constraints to 10.² In the case of CLUST, we set the number of clusters to 10, which we determine using the k-means algorithm. The entire code library, training and testing splits, and evaluation results will be made available online upon acceptance of the paper.

4 RESULTS AND DISCUSSION

We first summarize the results of the overall evaluation, and then discuss several aspects of these results in detail.

The overall cumulative match characteristic (CMC) curves for two representative single- and multi-shot datasets are shown in Figure 5. The CMC curve is a plot of the re-identification rate at rank-k. The individual performance of each algorithm combination on all datasets as well as complete CMC curves can be found in the supplementary material. As can be seen from the CMC curves, the “spread” in the performance of the algorithms for each dataset is huge, indicating the progress made by the re-id community over the past decade. However, on most datasets, the performance is still far from the point where we would consider re-id to be a solved problem.

2. This is set to 1 for kPCCA and rPCCA on Market1501 due to system memory issues.

Datasets	Best Combination	1	5	10
VIPeR	GOG-XQDA	41.1	71.1	82.1
GRID	GOG-XQDA	21.5	38.1	49.4
3DPeS	GOG-kMFA _l	51.0	77.0	86.9
CUHK03	GOG-kLFDA _{exp}	62.1	88.7	94.2
HDA+	AlexNet-FDA	81.6	82.4	83.2
Market1501	GOG-kLFDA _l	58.6	79.4	85.7
Airport	GOG-XQDA	34.8	58.2	69.0
PRID	GOG-KISSME-SRID	91.5	97.8	98.8
V47	GOG-KISSME-SRID	100.0	100.0	100.0
CAVIAR	GOG-KISSME-RNP	55.6	79.6	95.6
WARD-12	GOG-KISSME-SRID	99.7	100.0	100.0
WARD-13	GOG-KISSME-SRID	96.0	98.6	99.1
SAIVT-38	GOG-KISSME-SRID	96.5	100.0	100.0
SAIVT-58	GOG-KISSME-RNP	72.6	89.9	93.0
RAiD-12	GOG-KISSME-SRID	100.0	100.0	100.0
RAiD-13	GOG-KISSME-SRID	91.9	94.8	96.2
RAiD-14	GOG-KISSME-SRID	95.7	96.2	99.5
iLIDSVID	GOG-KISSME-SRID	75.7	90.1	93.6

TABLE 3

Top performing algorithmic combinations on each dataset. Read as feature-metric for single-shot and feature-metric-ranking for multi-shot.

In Table 3, we summarize the overall CMC curves by reporting the algorithm combination that achieved the best performance on each dataset as measured by the rank-1 performance. We note that GOG [44] performs the best among the 8 evaluated feature extraction algorithms, with it being a part of the best performing algorithm combination in 6 of the 7 single-shot and all the 11 multi-shot datasets. In the case of single-shot evaluation, while XQDA [18] is part of the best performing algorithm combination in 3 of the 7 single-shot datasets, across all the single-shot datasets, we note that GOG-kLFDA_l gives a mean rank-1 performance of 47.5%, marginally higher than GOG-XQDA, which gives 45.3%. This suggests that kLFDA_l is the best performing metric learning algorithm. In the case of multi-shot evaluation, the combination of KISSME [7] as the metric learning algorithm and SRID [52] as the multi-shot ranking algorithm is the best performing algorithm combination, with it resulting in the best performance on 9 of the 11

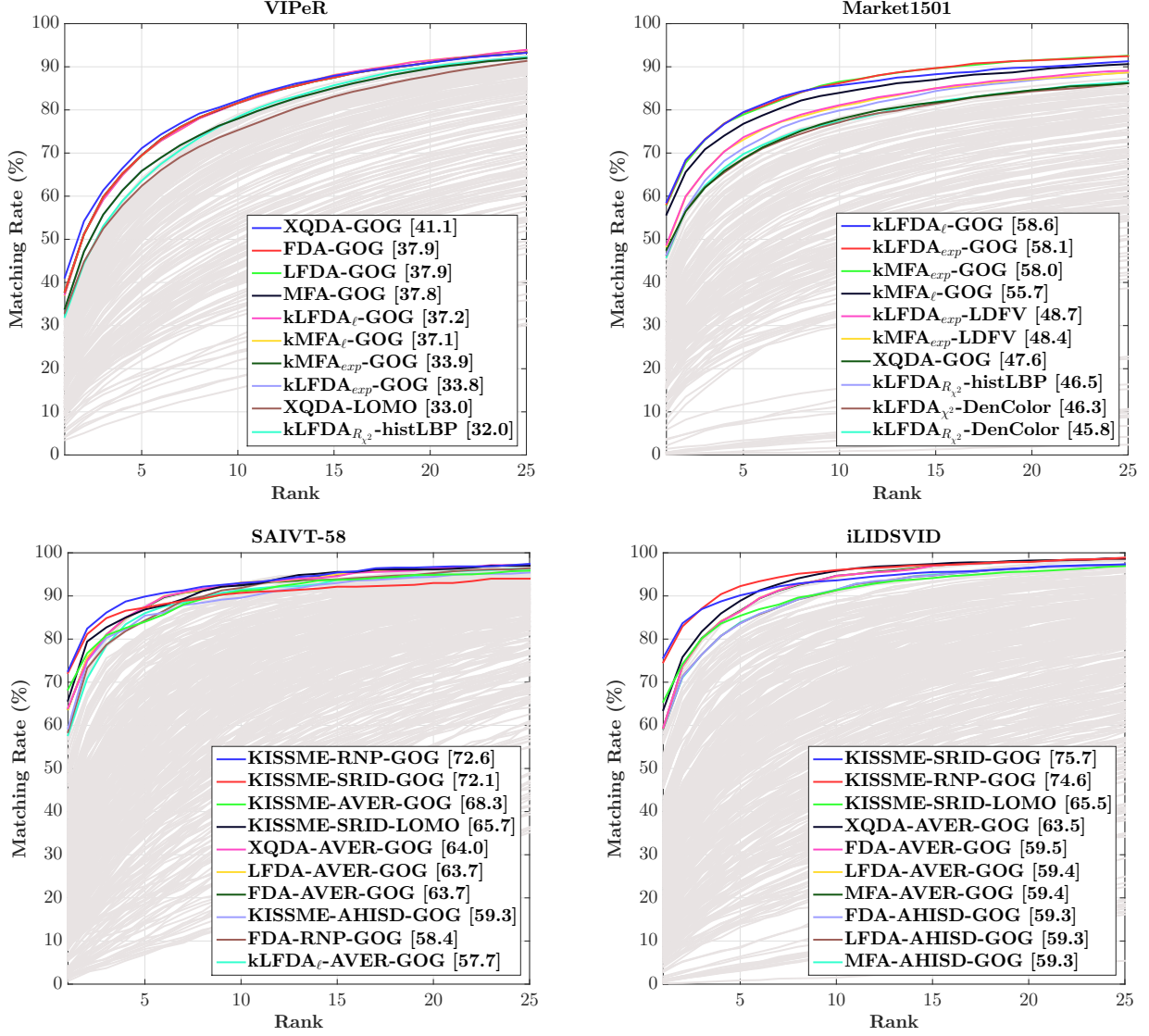


Fig. 5. CMC curves for the single-shot datasets VIPeR and Market1501, and the multi-shot datasets SAIVT-58 and iLIDSVID. The algorithmic combinations with the ten best rank-1 performances (indicated in the legend) are shown in color and all the others are shown in gray. CMC curves for all other datasets can be found in the supplementary material.

datasets.

In general, we observe that the algorithms give better performance on multi-shot datasets than on single-shot datasets. While this may be attributed to multi-shot datasets having more information in terms of the number of images per person, it is important to note that the single-shot datasets considered here generally have a significantly higher number of people in the gallery. It is quite natural to expect re-id performance to go down as the number of gallery people increases because we are now searching for the person of interest in a much larger candidate set.

4.1 Single shot analysis: features and metric learning

Single-shot re-id involves two key aspects: features and metric learning. In this section, we isolate the impact of the best performing algorithm in these two areas. First, we note that GOG is the best performing feature extraction algorithm in our evaluation. To corroborate this observation,

we study the impact of GOG in comparison with other feature extraction algorithms both in the presence as well as the absence of any metric learning. In the first experiment, we use the baseline Euclidean distance to rank gallery candidates in the originally computed feature space, which can be regarded as an unsupervised method. As can be noted from the results shown in Figure 6(a), GOG gives the best performance on 3 of the 7 datasets and very close performance to best performing algorithms on the other datasets.

Next, we study how GOG performs in comparison with other features in the presence of metric learning. In this experiment, we fix kLFDA_ℓ as our metric learning algorithm and rank gallery candidates using all the 7 evaluated feature extraction algorithms. The rank-1³ results for this

3. Complete CMC curves can be found in the supplementary material.

experiment are shown in Figure 6(b). As can be noted from the graph, GOG gives the best performance on 6 of the 7 datasets, with the exception of HDA where AlexNet gives the best performance.

These experiments clearly show that GOG is indeed the best performing feature extraction algorithm. This is because color and texture are the most descriptive aspects of a person image and GOG describes the global color and texture distributions using a local Gaussian distributions of pixel-level features. Another critical reason for the success of GOG is the hierarchical modeling of local color and texture structures. This is a critical step because typically a person’s clothes consists of local parts, each of which has certain local properties.

Next, we analyze the performance of different metric learning algorithms⁴, in the context of GOG, the best performing feature extraction algorithm. In this experiment, we fix GOG as the feature extraction algorithm and study how different metric learning algorithms perform. The results of this experiment are shown in Figure 6(c), from which we can note that XQDA gives the best performance on VIPeR, GRID, and Airport, whereas kLFDA give the best performance on CUKH03 and Market1501, suggesting both these metric learning algorithms are quite strong, further corroborating what we observed in Table 3. If we were to consider the mean rank-1 performance across all the 7 datasets, kLFDA₁ gives 47% whereas XQDA gives 45.3%, suggesting that kLFDA_{exp} is the best performing metric learning algorithm. Finally, to conclude the empirical study in this section, we analyze the choice of the kernel in kLFDA. In this experiment, we use GOG with the exponential and linear kernels, the results of which are shown in Figure 6(d). As can be seen from the bar graph, we do not observe any consistent trends across the 7 datasets. We note that kLFDA₁ gives a mean rank-1 performance of 47.5% whereas that of kLFDA_{exp} is 47%. While the difference is not statistically significant, this does illustrate that the extra non-linearity offered by the exponential kernel does not seem to help.

From the above discussion, we can infer the following: both kLFDA and XQDA emerge as strong metric learning algorithms, with GOG-kLFDA₁ giving a higher mean rank-1 performance across all the datasets when compared to GOG-XQDA. It is interesting to note that both kLFDA and XQDA involve solving generalized eigenvalue decomposition problems, similar to traditional Fisher discriminant analysis. This suggests that the approach of formulating discriminant objective functions in terms of data scatter matrices is most suitable to the re-id problem.

4.2 Multi-shot analysis: features, metric learning, and ranking

Multi-shot re-id involves three aspects: features, metric learning, and ranking. As noted previously, GOG, KISSME, and SRID emerged as the best performing algorithmic combination. On all the datasets, as expected, a custom ranking algorithm resulted in the best performance, with SRID performing the best on 9 of these 11 datasets. In this section, we provide further empirical results analyzing the

impact of using a multi-shot ranking algorithm. To this end, we fix GOG as the feature extraction scheme.

First, we evaluate the impact of using a multi-shot ranking algorithm instead of AVER. Here, we compare the performance of GOG-KISSME-AVER and GOG-KISSME-SRID. The results are shown in Figure 7(b). As can be noted from the graph, with the exception of CAVIAR, SRID generally gives superior performance when compared to AVER. This, and our observations from Table 3, suggest that using a multi-shot ranking algorithm that exploits the inherent structure of the data instead of naive feature averaging will give better performance. Furthermore, we also note that a multi-shot ranking algorithm in itself is not sufficient to give good performance because that would be purely an unsupervised approach. Combining a metric learning algorithm with the ranking technique adds a layer of supervision to the overall framework and will provide a significant performance boost.

Next, we analyze the performance of the feature extraction and metric learning algorithms and compare the observed trends with those in the single-shot case. In the feature extraction case, we fix SRID as the ranking algorithm and KISSME as the metric learning algorithm. The rank-1 results for this experiment are shown in Figure 7(b)-(c). We see very clear trends in this case, with GOG giving the best results on all the 11 datasets. These results are not surprising given the strong performance of GOG in the single-shot case. In the metric learning case, we fix SRID as the ranking algorithm and GOG as the feature extraction algorithm, with Figure 7(d) showing the rank-1 results. We see very clear trends in this case as well, with KISSME giving the best results across all datasets.

4.3 Additional observations: PCA, features, and attributes

In this section, we report additional empirical observations. Most contemporary feature extraction schemes produce high-dimensional data, introducing significant computational overhead. To this end, we analyze the impact of an unsupervised dimensionality reduction scheme, principal component analysis (PCA). We fix GOG as the feature extraction scheme and perform experiments with and without PCA. We set the dimension of the PCA-reduced space to 100. The results are shown in the first two bars (pink and yellow) in Figure 8(a) and 8(b). The results without PCA are better than those with PCA on all the 18 datasets. This observation is not surprising given that PCA can result in the undesirable removal of the most discriminative features.

Next, we analyze the impact of the number-of-strips parameter in the best feature extraction algorithm, GOG. To this end, we perform experiments with 6, 9, and 15 horizontal strips in GOG, with Euclidean distance as the metric in the single-shot case and Euclidean distance as the metric and AVER as the ranking strategy in the multi-shot case. The rank-1 results are shown in bars 2–5 in Figure 8(a) and 8(b). While it is reasonable to expect superior performance as we increase the number of strips, thereby increasing the feature space dimensionality, it is important to note that in this process, we may have fewer samples to estimate the Gaussians in each strip and also increase the amount of

4. A discussion on the training time of these algorithms is provided in the supplementary material.

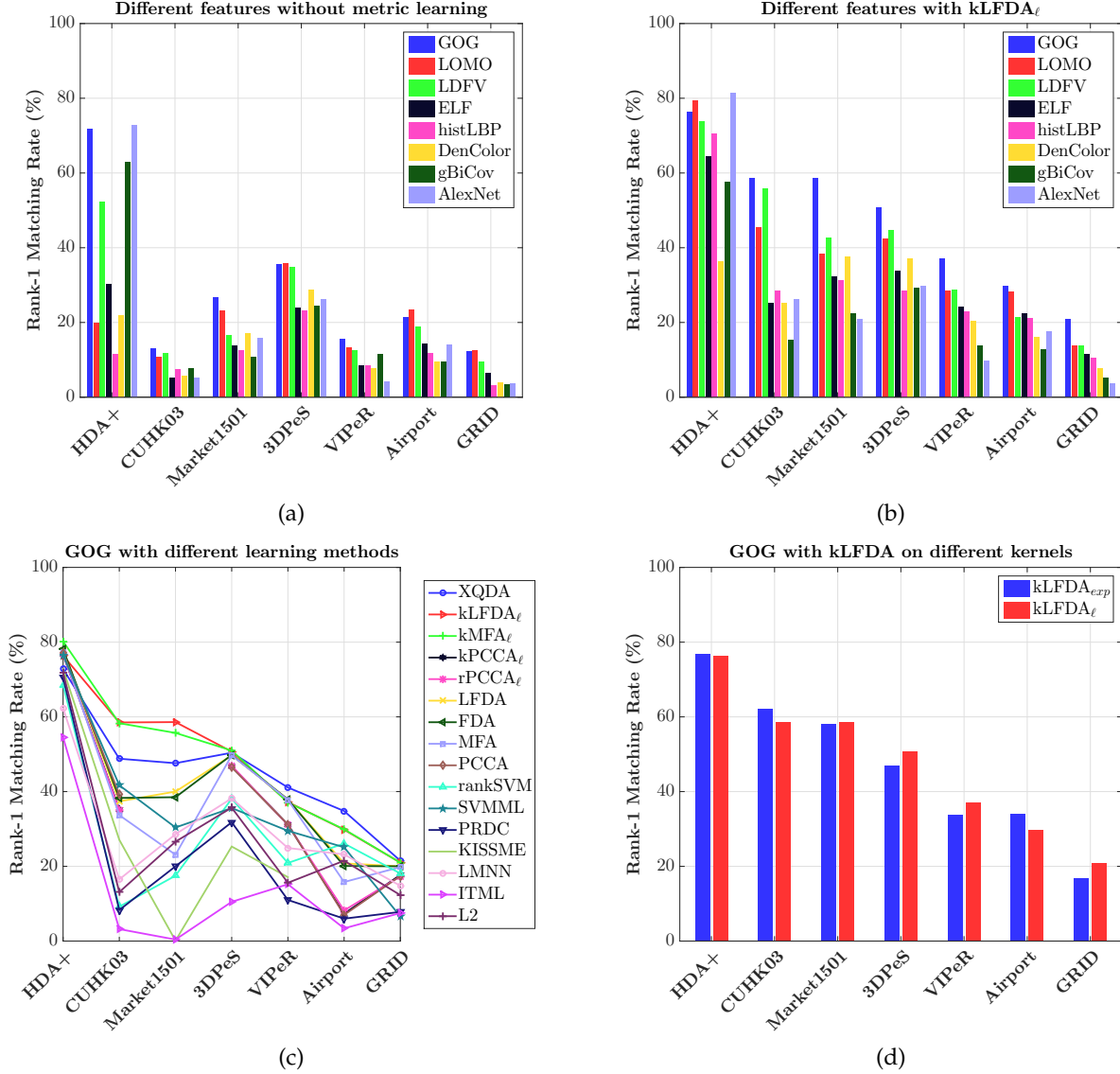


Fig. 6. Rank-1 results for single shot datasets illustrating the impact of GOG and kLFDA, the best performing feature extraction and metric learning algorithms.

background/noise/non-informative content in the feature descriptor. This indeed seems to be the case on many multi-shot datasets, with performance generally decreasing (except for CAVIAR and RAiD-12), albeit marginally, as we increase the number of strips, as can be noted from bars 2–4. In the single-shot case, there does not seem to be any significant performance variations as we increase the number of strips. Given the computational complexity involved in working with higher dimensional feature spaces due to increased number of strips, these results suggest that 6 strips, which is in fact the widely used number in the re-id community, seems to be a reasonable choice, giving better or close performance to the other choices in most cases.

Finally, we analyze the performance of the different feature extraction schemes with respect to the different attributes used to characterize datasets in Table 2. The goal of this experiment is to study which features are good in certain scenarios. To this end, we use Euclidean distance as

the metric, and in the multi-shot case, AVER as the ranking algorithm, and report the mean rank-1 performance on all datasets for each attribute group. The results obtained are shown in Figure 8(c) and 8(d) for the single-shot and multi-shot cases, respectively. We observe the following trends from the results. In the single-shot case, GOG resulted in the best performance on VV, IV, and RES, whereas LOMO gave the best performance on BC and very close performance on OCC. We also note that LDFV gives strong performance, behind GOG, on VV and IV. In the multi-shot case, GOG resulted in the best performance on all the five attributes. These results suggest that modeling local pixel distributions, as done in all these three methods, is an extremely important step in describing person images. Additionally, since viewpoint invariance is an extremely important attribute for any re-id descriptor, these results suggest that incorporating some sort of local region covariance information, as done in both GOG and LOMO is critical to achieve viewpoint

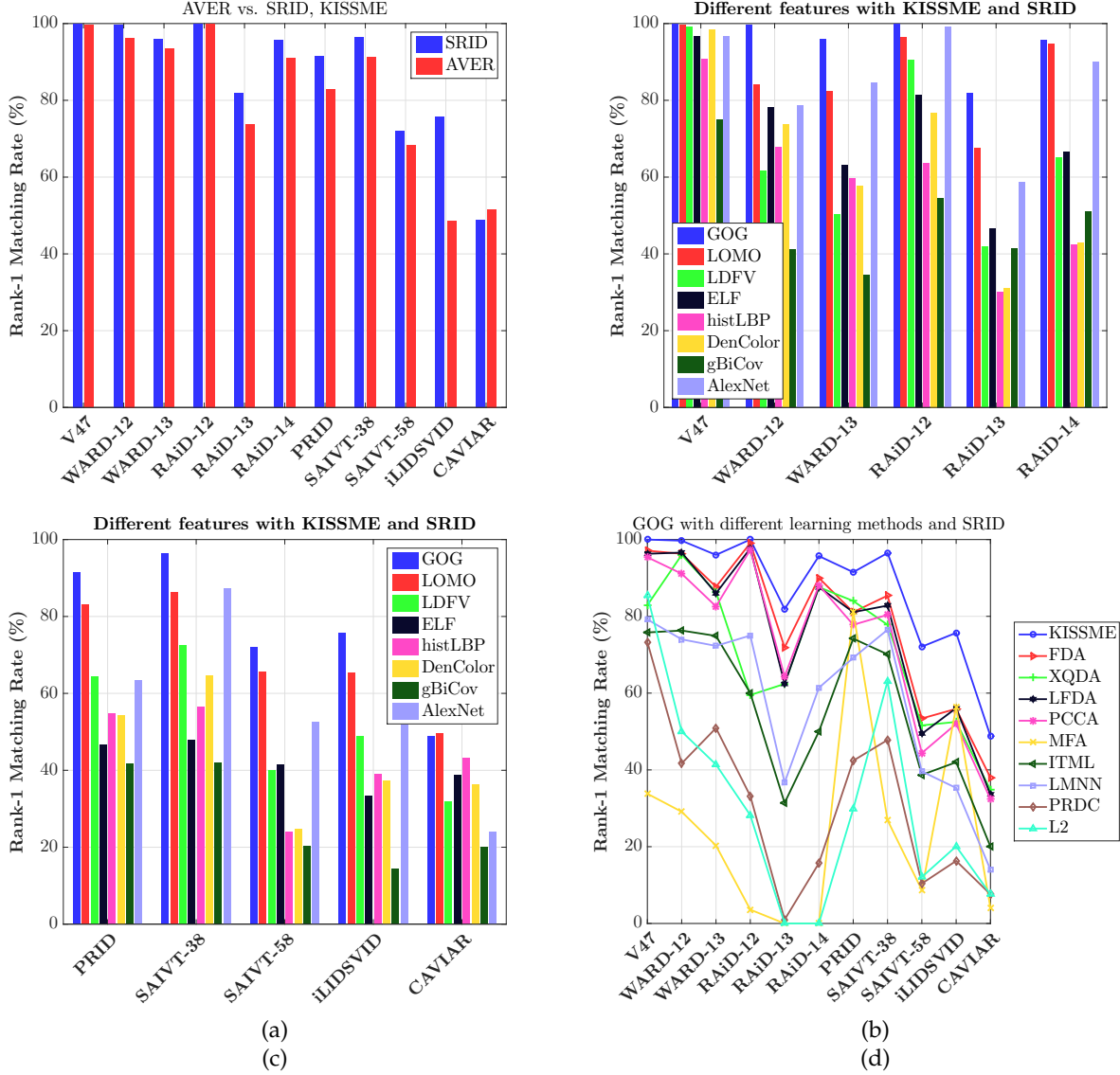


Fig. 7. (a): Rank-1 performance on multi-shot datasets, illustrating the impact of the best performing multi-shot ranking algorithm, SRID over AVER, naive feature averaging. (b)-(d) Rank-1 performance on multi-shot datasets comparing various feature extraction and metric learning algorithms with SRID as the ranking algorithm.

invariant descriptors. Furthermore, we note that LOMO results in strong performance on BC and OCC in the single-shot case. This is primarily due to the presence of SILTP features [42] in the descriptor, which adds a certain degree of robustness against noise and other corruptions in images. However, we do not see a similar trend in the multi-shot case, where the bar graph is dominated by the performance of GOG on iLIDSVID, which, in addition to background clutter and occlusions, has large viewpoint and illumination variations. While deep learning descriptors have performed exceedingly well on tasks such as object recognition, in the context of re-id as seen from these graphs, it's performance is not as high as some of the other methods evaluated. This is partly due to the lack of diversity in the re-id datasets. Some suggestions for mitigating this issue are provided in Section 5. However, we do note that AlexNet gives strong performance behind GOG on the IV attribute, and is in fact

better than LOMO in the multi-shot case. While LOMO uses the Retinex transform to obtain illumination invariance, the strong relative performance of AlexNet suggests that learning CNN-based descriptors may be better suited to such scenarios. Incorporating some of the recent advances in architecture design [76], [77] can potentially lead to further performance gains.

4.4 Impact of datasets on re-id research

Datasets play an extremely important role in validating a newly proposed algorithm. From Table 3, we note that the V47 dataset can be considered to be solved, with the best performing algorithm achieving a rank-1 performance of 100%. However, the performance on other multi-shot datasets is still far from ideal. These datasets therefore present opportunities and challenges to develop better algorithms. A similar argument holds for the single-shot case.

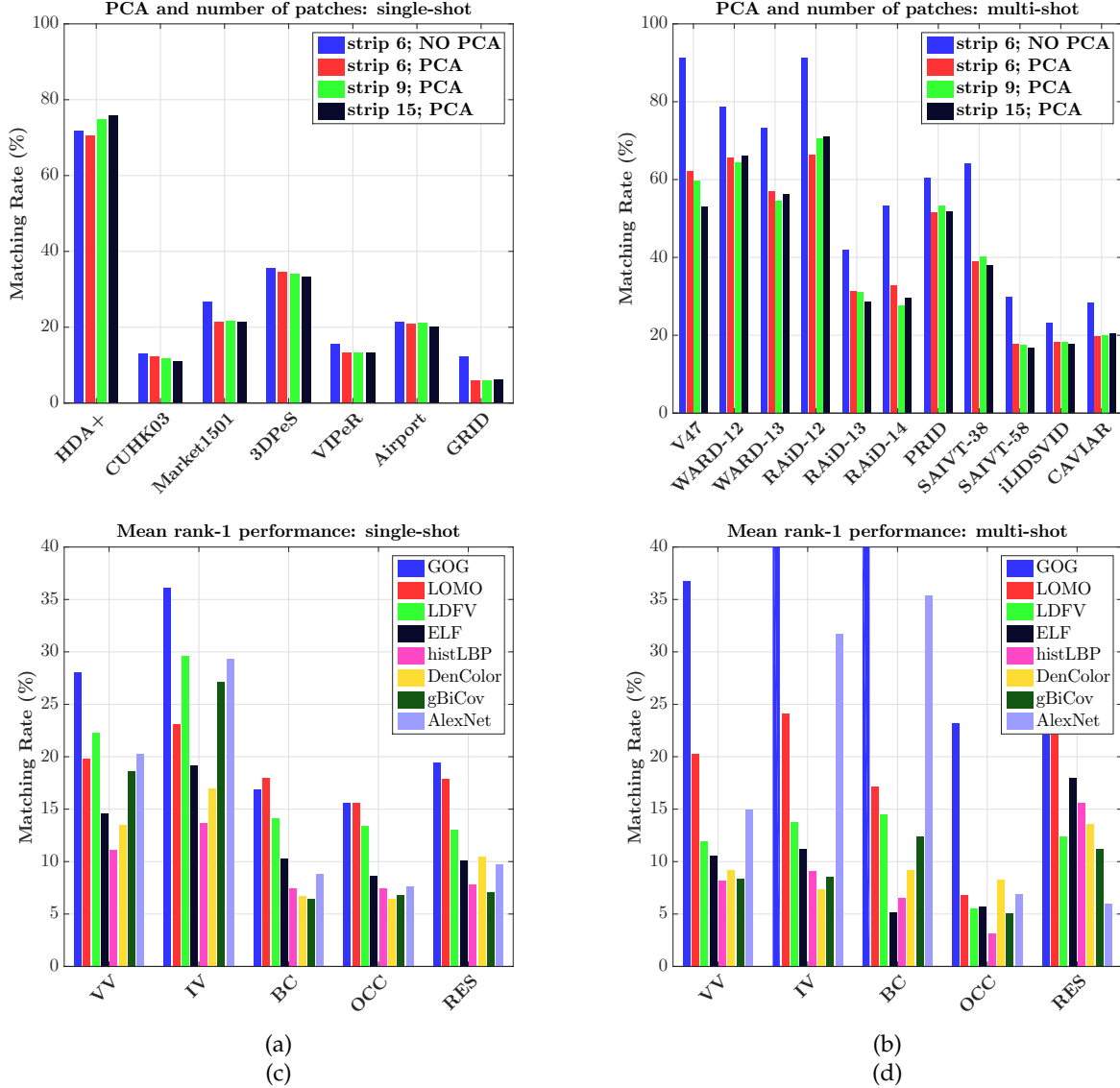


Fig. 8. Rank-1 performance on (a) single-shot and (b) multi-shot datasets illustrating the impact of PCA and number of strips. (c)-(d) Mean rank-1 performance across all single- and multi-shot datasets with respect to various attributes and features.

The performance on VIPeR is still very low despite it being the most popular dataset in the re-id community. The performance on GRID is the lowest (21.5% at rank-1) and this is in part due to the presence of a large number of distractor people in the gallery. The newly proposed Airport dataset has the next lowest performance (34.8% at rank-1). This is due to the presence of a large distractor set as well as the significant real-world challenges described in Section 3.1. These observations suggest that as the number of distractor people in the gallery increases, the performance of a re-id algorithm goes down. This is not surprising since now, we have a larger set of people to compare the probe person against, leading to more avenues for the re-id algorithm to fail. The overall performance on all the datasets can be improved by designing better algorithms, some suggestions for which are discussed in Section 5.

Our categorization of datasets according to their attributes also provides a useful reference point for construct-

ing new datasets. We note that none of the 14 datasets have all 6 attributes. While MARS [78], a recently proposed dataset, has a large number of images, constructing datasets that are of the size of ImageNet, in terms of the number of people and positive examples, will assist in the application of some of the recent algorithmic advances in feature learning using CNNs [76], [77]. We believe focusing on these aspects while collecting new datasets would help accelerate progress in person re-id.

5 CONCLUSIONS AND FUTURE RESEARCH DIRECTIONS

In this work, we presented a new, large-scale dataset that accurately reflects the real-world person re-identification problem. We conducted large-scale experiments on the new dataset as well as 13 existing datasets, producing a systematic re-id evaluation and benchmark that we believe will be extremely valuable to the re-id community.

Based on our results and analysis, we posit important research directions that could be potential candidates for developing better algorithms to improve re-id performance.

- We noted that GOG resulted in the best performance among all the evaluated feature extraction methods. A primary reason for its success is the hierarchical modeling of local pixel distributions in terms of local means and covariances. While such local covariance-based features have previously been constructed [79], learning covariant feature detectors in an end-to-end fashion following recent advances [80] can lead to further performance gains.
- In feature extraction, a typical approach is to concatenate local features from all image patches or partitions. However, LOMO, with its local max-pooling, resulted in strong performance when compared to features such as ELF that employ such local feature concatenation. Clearly, such local max-pooling, which is already widely used in convolutional neural networks [40], is important and can be a useful tool in the context of re-id. Furthermore, as done in GOG, modeling local color and texture features is critical for a robust descriptor. This suggests developing end-to-end texture representation schemes [81], [82] could be a promising research direction in developing descriptors for re-id.
- In metric learning, we noted that both kLFDA and XQDA gave strong performance. Since both these methods formulate Fisher-type criteria to maximize the class likelihood of individual samples, a natural extension would be to learn, in an end-to-end fashion, CNN representations that minimize within-class variance while maximizing between-class variance [83].
- In the context of multi-shot data, developing feature extraction or learning schemes that specifically exploit the temporal structure of the data [28], [84], [85] will be a natural extension of the existing methods that use only spatial image information.
- In multi-shot ranking, we demonstrated that using a custom ranking method gives much better performance when compared to using the feature averaging scheme. In practical re-id applications with multi-shot data, an image sequence of a person will typically undergo several variations such as background clutter, occlusion and illumination variations, and developing custom multi-shot ranking algorithms that take all this data variance into account will give better performance. Another promising future research direction in this context would be to integrate multi-shot ranking with metric learning. While most existing methods treat these two topics separately, developing a unified metric learning and multi-shot ranking framework that exploits the several aspects of multi-shot data can potentially lead to further performance gains. For instance, borrowing ideas from research in spatio-temporal feature learning [86], [87] would be a natural next step in developing such unified algorithms.
- As noted in Section 4.4, most existing datasets are

quite small and constrained to learn an effective generic feature space using convolutional neural networks, due to which we have to resort to finetuning existing models. Furthermore, most datasets are typically captured under fixed conditions, such as specific indoor/outdoor locations and backgrounds. Constructing datasets that contain a much larger number of people than existing datasets as well as capturing images under an eclectic mix of conditions will rapidly accelerate progress in learning generic feature descriptors for re-id.

ACKNOWLEDGEMENTS

This material is based upon work supported by the U.S. Department of Homeland Security under Award Number 2013- ST-061-ED0001. The views and conclusions contained in this document are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of the U.S. Department of Homeland Security. Thanks to Michael Young, Jim Spriggs, and Don Kemer for supplying the airport video data.

REFERENCES

- [1] Y. Li, Z. Wu, S. Karanam, and R. Radke, "Real-world re-identification in an airport camera network," in *ICDSC*, 2014.
- [2] S. Gong, M. Cristani, S. Yan, and C. C. Loy, *Person re-identification*. Springer, 2014, vol. 1.
- [3] O. Camps, M. Gou, T. Hebble, S. Karanam, O. Lehmann, Y. Li, R. Radke, Z. Wu, and F. Xiong, "From the lab to the real world: Re-identification in an airport camera network," *IEEE Transactions on Circuits and Systems for Video Technology (T-CSVT)*, vol. PP, no. 99, 2016.
- [4] D. Gray and H. Tao, "Viewpoint invariant pedestrian recognition with an ensemble of localized features," in *ECCV*, 2008.
- [5] B. Prosser, W.-S. Zheng, S. Gong, and T. Xiang, "Person re-identification by support vector ranking," in *BMVC*, 2010.
- [6] W.-S. Zheng, S. Gong, and T. Xiang, "Person re-identification by probabilistic relative distance comparison," in *CVPR*, 2011.
- [7] M. Koestinger, M. Hirzer, P. Wohlhart, P. M. Roth, and H. Bischof, "Large scale metric learning from equivalence constraints," in *CVPR*, 2012.
- [8] A. Mignon and F. Jurie, "PCCA: A new approach for distance learning from sparse pairwise constraints," in *CVPR*, 2012.
- [9] R. Zhao, W. Ouyang, and X. Wang, "Unsupervised salience learning for person re-identification," in *CVPR*, 2013.
- [10] L. Bazzani, M. Cristani, and V. Murino, "Symmetry-driven accumulation of local features for human characterization and re-identification," *CVIU*, vol. 117, no. 2, pp. 130–144, 2013.
- [11] S. Pedagadi, J. Orwell, S. Velastin, and B. Boghossian, "Local fisher discriminant analysis for pedestrian re-identification," in *CVPR*, 2013.
- [12] L. An, M. Kafai, S. Yang, and B. Bhanu, "Reference-based person re-identification," in *AVSS*. IEEE, 2013, pp. 244–249.
- [13] F. Xiong, M. Gou, O. Camps, and M. Sznaier, "Person re-identification using kernel-based metric learning methods," in *ECCV*, 2014.
- [14] R. Zhao, W. Ouyang, and X. Wang, "Learning mid-level filters for person re-identification," in *CVPR*, 2014.
- [15] Z. Wu, Y. Li, and R. Radke, "Viewpoint invariant human re-identification in camera networks using pose priors and subject-discriminative features," *T-PAMI*, vol. 37, no. 5, pp. 1095–1108, 2015.
- [16] S. Paisitkriangkrai, C. Shen, and A. van den Hengel, "Learning to rank in person re-identification with metric ensembles," in *CVPR*, 2015.
- [17] S. Karanam, Y. Li, and R. J. Radke, "Person re-identification with discriminatively trained viewpoint invariant dictionaries," in *ICCV*, 2015.

- [18] S. Liao, Y. Hu, X. Zhu, and S. Z. Li, "Person re-identification by local maximal occurrence representation and metric learning," in *CVPR*, 2015.
- [19] W.-S. Zheng, X. Li, T. Xiang, S. Liao, J. Lai, and S. Gong, "Partial person re-identification," in *ICCV*, 2015.
- [20] X. Li, W.-S. Zheng, X. Wang, T. Xiang, and S. Gong, "Multi-scale learning for low-resolution person re-identification," in *ICCV*, 2015.
- [21] S. Messelodi and C. M. Modena, "Boosting fisher vector based scoring functions for person re-identification," *Image and Vision Computing (IVC)*, vol. 44, pp. 44–58, 2015.
- [22] D. Chen, Z. Yuan, G. Hua, N. Zheng, and J. Wang, "Similarity learning on an explicit polynomial kernel feature map for person re-identification," in *CVPR*, 2015, pp. 1565–1573.
- [23] M. Gou, X. Zhang, A. Rates-Borras, S. Asghari-Esfeden, M. Sznaiier, and O. Camps, "Person re-identification in appearance impaired scenarios," in *BMVC*, 2016.
- [24] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino, "Custom pictorial structures for re-identification," in *BMVC*, 2011.
- [25] M. Hirzer, C. Beleznaï, P. M. Roth, and H. Bischof, "Person re-identification by descriptive and discriminative classification," in *Image Analysis*, 2011.
- [26] D. Baltieri, R. Vezzani, and R. Cucchiara, "3DPeS: 3d people dataset for surveillance and forensics," in *Proceedings of the 2011 joint ACM workshop on Human gesture and behavior understanding*, 2011.
- [27] N. Martinel, C. Micheloni, and C. Piciarelli, "Distributed signature fusion for person re-identification," in *ICDSC*, 2012.
- [28] T. Wang, S. Gong, X. Zhu, and S. Wang, "Person re-identification by video ranking," in *ECCV*, 2014.
- [29] A. Das, A. Chakraborty, and A. K. Roy-Chowdhury, "Consistent re-identification in a camera network," in *ECCV*, 2014.
- [30] C. C. Loy, C. Liu, and S. Gong, "Person re-identification by manifold ranking," in *ICIP*, 2013.
- [31] D. Scharstein and R. Szeliski, "A taxonomy and evaluation of dense two-frame stereo correspondence algorithms," *IJCV*, vol. 47, no. 1–3, pp. 7–42, 2002.
- [32] S. Baker, D. Scharstein, J. Lewis, S. Roth, M. J. Black, and R. Szeliski, "A database and evaluation methodology for optical flow," *IJCV*, vol. 92, no. 1, pp. 1–31, 2011.
- [33] C. Schmid, "Constructing models for content-based image retrieval," in *CVPR*, 2001.
- [34] I. Fogel and D. Sagi, "Gabor filters as texture discriminator," *Biological cybernetics*, vol. 61, no. 2, pp. 103–113, 1989.
- [35] B. Ma, Y. Su, and F. Jurie, "Local descriptors encoded by fisher vectors for person re-identification," in *ECCV Workshops*, 2012.
- [36] J. Sánchez, F. Perronnin, T. Mensink, and J. Verbeek, "Image classification with the fisher vector: Theory and practice," *IJCV*, vol. 105, no. 3, pp. 222–245, 2013.
- [37] B. Ma, Y. Su, and F. Jurie, "Covariance descriptor based on bio-inspired features for person re-identification and face verification," *Image and Vision Computing (IVC)*, vol. 32, no. 6, pp. 379–390, 2014.
- [38] M. Riesenhuber and T. Poggio, "Hierarchical models of object recognition in cortex," *Nature neuroscience*, vol. 2, no. 11, pp. 1019–1025, 1999.
- [39] O. Tuzel, F. Porikli, and P. Meer, "Pedestrian detection via classification on riemannian manifolds," *T-PAMI*, vol. 30, no. 10, pp. 1713–1727, 2008.
- [40] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, 2012.
- [41] T. Ojala, M. Pietikäinen, and D. Harwood, "A comparative study of texture measures with classification based on featured distributions," *Pattern recognition*, vol. 29, no. 1, pp. 51–59, 1996.
- [42] S. Liao, G. Zhao, V. Kellokumpu, M. Pietikäinen, and S. Z. Li, "Modeling pixel process with scale invariant local patterns for background subtraction in complex scenes," in *CVPR*, 2010.
- [43] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *T-IP*, vol. 6, no. 7, pp. 965–976, 1997.
- [44] T. Matsukawa, T. Okabe, E. Suzuki, and Y. Sato, "Hierarchical gaussian descriptor for person re-identification," in *CVPR*, 2016.
- [45] R. A. Fisher, "The use of multiple measurements in taxonomic problems," *Annals of eugenics (AE)*, vol. 7, no. 2, pp. 179–188, 1936.
- [46] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *ICML*, 2007.
- [47] S. Yan, D. Xu, B. Zhang, H.-J. Zhang, Q. Yang, and S. Lin, "Graph embedding and extensions: a general framework for dimensionality reduction," *T-PAMI*, vol. 29, no. 1, pp. 40–51, 2007.
- [48] Z. Li, S. Chang, F. Liang, T. S. Huang, L. Cao, and J. R. Smith, "Learning locally-adaptive decision functions for person verification," in *CVPR*, 2013.
- [49] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *JMLR*, vol. 10, pp. 207–244, 2009.
- [50] H. Cevikalp and B. Triggs, "Face recognition based on image sets," in *CVPR*, 2010.
- [51] M. Yang, P. Zhu, L. Van Gool, and L. Zhang, "Face recognition based on regularized nearest points between image sets," in *FG*, 2013.
- [52] S. Karanam, Y. Li, and R. Radke, "Sparse re-id: Block sparsity for person re-identification," in *CVPR Workshops*, 2015.
- [53] C. Liu, C. C. Loy, S. Gong, and G. Wang, "POP: Person re-identification post-rank optimisation," in *ICCV*, 2013.
- [54] J. Garcia, N. Martinel, C. Micheloni, and A. Gardel, "Person re-identification ranking optimisation by discriminant context information analysis," in *ICCV*, 2015.
- [55] S. Bak, E. Corvee, F. Bremond, and M. Thonnat, "Multiple-shot human re-identification by mean riemannian covariance grid," in *AVSS*, 2011.
- [56] C. Su, F. Yang, S. Zhang, Q. Tian, L. S. Davis, and W. Gao, "Multi-task learning with low rank attribute embedding for person re-identification," in *ICCV*, 2015.
- [57] R. Layne, T. M. Hospedales, and S. Gong, "Person re-identification by attributes," in *BMVC*, 2012.
- [58] Z. Shi, T. M. Hospedales, and T. Xiang, "Transferring a semantic representation for person re-identification and search," in *CVPR*, 2015.
- [59] M. Eisenbach, A. Kolarow, A. Vorndran, J. Niebling, and H.-M. Gross, "Evaluation of multi feature fusion at score-level for appearance-based person re-identification," in *Int. Joint Conf. Neural Networks (IJCNN)*, 2015.
- [60] L. Zheng, S. Wang, L. Tian, F. He, Z. Liu, and Q. Tian, "Query-adaptive late fusion for image search and person re-identification," in *CVPR*, 2015.
- [61] N. Martinel, C. Micheloni, and G. L. Foresti, "A pool of multiple person re-identification experts," *Pattern Recognition Lett. (PRL)*, vol. 71, pp. 23–30, 2016.
- [62] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification: Past, present and future," *arXiv preprint arXiv:1610.02984*, 2016.
- [63] R. Satta, "Appearance descriptors for person re-identification: a comprehensive review," *arXiv preprint arXiv:1307.5748*, 2013.
- [64] R. Vezzani, D. Baltieri, and R. Cucchiara, "People reidentification in surveillance and forensics: A survey," *ACM Computing Surveys (CSUR)*, vol. 46, no. 2, p. 29, 2013.
- [65] A. Bedagkar-Gala and S. K. Shah, "A survey of approaches and trends in person re-identification," *Image and Vision Computing*, vol. 32, no. 4, pp. 270–286, 2014.
- [66] P. Dollár, R. Appel, S. Belongie, and P. Perona, "Fast feature pyramids for object detection," *T-PAMI*, vol. 36, no. 8, pp. 1532–1545, 2014.
- [67] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, "Object detection with discriminatively trained part-based models," *T-PAMI*, vol. 32, no. 9, pp. 1627–1645, 2010.
- [68] C. C. Loy, T. Xiang, and S. Gong, "Time-delayed correlation analysis for multi-camera activity understanding," *IJCV*, vol. 90, no. 1, pp. 106–129, 2010.
- [69] S. Wang, M. Lewandowski, J. Annesley, and J. Orwell, "Re-identification of pedestrians with variable occlusion and scale," in *ICCV Workshops*, 2011.
- [70] A. Bialkowski, S. Denman, S. Sridharan, C. Fookes, and P. Lucey, "A database for person re-identification in multi-camera surveillance networks," in *DICTA*, 2012.
- [71] W. Li, R. Zhao, T. Xiao, and X. Wang, "DeepReId: Deep filter pairing neural network for person re-identification," in *CVPR*, 2014.
- [72] D. Figueira, M. Taiana, A. Nambiar, J. Nascimento, and A. Bernardino, "The HDA+ data set for research on fully automated re-identification systems," in *ECCV Workshops*, 2014.
- [73] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, and Q. Tian, "Scalable person re-identification: A benchmark," in *ICCV*, 2015.

- [74] E. Rosten, R. Porter, and T. Drummond, "Faster and better: A machine learning approach to corner detection," *T-PAMI*, vol. 32, no. 1, pp. 105–119, 2010.
- [75] B. D. Lucas and T. Kanade, "An iterative image registration technique with an application to stereo vision," in *Imaging Understanding Workshop*, 1981.
- [76] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [77] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *arXiv preprint arXiv:1512.03385*, 2015.
- [78] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, and Q. Tian, "MARS: A video benchmark for large-scale person re-identification," in *ECCV*, 2016.
- [79] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification," in *ECCV*, 2006.
- [80] K. Lenc and A. Vedaldi, "Learning covariant feature detectors," *arXiv preprint arXiv:1605.01224*, 2016.
- [81] M. Cimpoi, S. Maji, I. Kokkinos, and A. Vedaldi, "Deep filter banks for texture recognition, description, and segmentation," *IJCV*, vol. 118, no. 1, pp. 65–94, 2016.
- [82] T.-Y. Lin and S. Maji, "Visualizing and understanding deep texture representations," in *CVPR*, 2016, pp. 2791–2799.
- [83] M. Dorfer, R. Kelz, and G. Widmer, "Deep linear discriminant analysis," in *ICLR*, 2016.
- [84] A. Klaser, M. Marszałek, and C. Schmid, "A spatio-temporal descriptor based on 3d-gradients," in *BMVC*, 2008.
- [85] K. Liu, B. Ma, W. Zhang, and R. Huang, "A spatio-temporal appearance representation for video-based pedestrian re-identification," in *ICCV*, 2015.
- [86] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3d convolutional networks," in *ICCV*, 2015.
- [87] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool, "Temporal segment networks: towards good practices for deep action recognition," in *ECCV*, 2016.



Angels Rates-Borras Angels Rates-Borras received her B.S. degree in Electrical Engineering from Universitat Politècnica de Catalunya, Barcelona, Spain. She is currently a Ph.D. student in the Department of Electrical and Computer Engineering at Northeastern University, Boston. Her research interests are mainly focused on person re-identification, scene understanding, video analysis and activity recognition.



Octavia Camps Octavia Camps received B.S. degrees in computer science in 1981 and in electrical engineering in 1984, from the Universidad de la República (Montevideo, Uruguay), and M.S. and Ph.D. degrees in electrical engineering in 1987 and 1992, from the University of Washington, respectively. Since 2006 she is a Professor in the Electrical and Computer Engineering Department at Northeastern University. From 1991 to 2006 she was a faculty member at the departments of Electrical Engineering and Computer Science and Engineering at The Pennsylvania State University. In 2000, she was a visiting faculty at the California Institute of Technology and at the University of Southern California and in 2013 she was a visiting faculty at the Computer Science Department at Boston University. Her main research interests include dynamics-based computer vision, image processing, and machine learning. She is a member of IEEE.



Srikrishna Karanam Srikrishna Karanam is a Ph.D. candidate in the Department of Electrical, Computer, and Systems Engineering at Rensselaer Polytechnic Institute (RPI), Troy, NY. He received the B.Tech. degree in Electronics and Communication Engineering from the National Institute of Technology Warangal in 2013 and the M.S. degree in Electrical Engineering from RPI in 2014. His research interests include computer vision, video processing, machine learning, and optimization.



Mengran Guo Mengran Guo is currently a Ph.D. candidate in the Department of Electrical and Computer Engineering at Northeastern University. He received an M.S. degree from the Pennsylvania State University and B.Eng. degree from Harbin Institute of Technology in China. His research interests are person re-identification and activity recognition.



Ziyang Wu Ziyang Wu received a Ph.D. degree in Computer and Systems Engineering from Rensselaer Polytechnic Institute in 2014. He has B.S. and M.S. degrees in Engineering from Beihang University. He joined Siemens Corporate Research as a Research Scientist in 2014. His current research interests include 3D object recognition and autonomous perception.



Richard J. Radke Richard J. Radke joined the Electrical, Computer, and Systems Engineering department at Rensselaer Polytechnic Institute in 2001, where he is now a Full Professor. He has B.A. and M.A. degrees in computational and applied mathematics from Rice University, and M.A. and Ph.D. degrees in electrical engineering from Princeton University. His current research interests involve computer vision problems related to human-scale, occupant-aware environments, such as person tracking and re-identification with cameras and range sensors. Dr. Radke is affiliated with the NSF Engineering Research Center for Lighting Enabled Service and Applications (LESA), the DHS Center of Excellence on Explosives Detection, Mitigation and Response (ALERT), and Rensselaer's Experimental Media and Performing Arts Center (EMPAC). He received an NSF CAREER award in March 2003 and was a member of the 2007 DARPA Computer Science Study Group. Dr. Radke is a Senior Member of the IEEE and a Senior Area Editor of *IEEE Transactions on Image Processing*. His textbook *Computer Vision for Visual Effects* was published by Cambridge University Press in 2012.