# Exploiting feature representations through similarity learning, post-ranking and ranking aggregation for person re-identification

Julio C. S. Jacques Junior[a,b,*], Xavier Baró[a,b], Sergio Escalera[c,b]

[a]*Faculty of Computer Science, Multimedia and Telecommunication - Universitat Oberta de Catalunya, Spain*
[b]*Computer Vision Center - Universitat Autònoma de Barcelona, Spain*
[c]*Department of Mathematics and Informatics - University of Barcelona, Spain*

## Abstract

Person re-identification has received special attention by the human analysis community in the last few years. To address the challenges in this field, many researchers have proposed different strategies, which basically exploit either cross-view invariant features or cross-view robust metrics. In this work, we propose to exploit a post-ranking approach and combine different feature representations through ranking aggregation. Spatial information, which potentially benefits the person matching, is represented using a 2D body model, from which color and texture information are extracted and combined. We also consider background/foreground information, automatically extracted via Deep Decompositional Network, and the usage of Convolutional Neural Network (CNN) features. To describe the matching between images we use the polynomial feature map, also taking into account local and global information. The Discriminant Context Information Analysis based post-ranking approach is used to improve initial ranking lists. Finally, the Stuart ranking aggregation method is employed to combine complementary ranking lists obtained from different feature representations. Experimental results demonstrated that we improve the state-of-the-art on VIPeR and PRID450s datasets, achieving 67.21% and 75.64% on top-1 rank recognition rate, respectively, as well as obtaining competitive results on CUHK01 dataset.

*Keywords:* person re-identification, similarity learning, feature fusion, post-ranking, ranking aggregation.

## 1. Introduction

Person re-identification is the task of assigning the same identifier to all instances of a particular individual captured in a series of images or videos, even after the occurrence of significant gaps over time or space. It has a wide range of applications, most of them focused on surveillance and forensic systems. Even though the proposed models and reported results in this field have considerably advanced in recent years [1, 2, 3], this task still presents open challenges, mainly due to the influence of numerous real-world factors such as illumination problems, occlusions, camera settings, as well as factors associated with the dynamics of the human being, like the great variety of appearance features, pose variations and strong visual similarity between different people. These difficulties are often compounded by low resolution images or poor quality video feeds with large amounts of unrelated information, making re-identification even harder.

As related in [4], given a query person image, in order to find the correct matches among a large set of candidate images captured by different cameras, two crucial problems have to be addressed. First, good image features are required to represent both the query and the gallery images. Second, suitable distance metrics are indispensable to determine whether a gallery image contains the same individual as the query image. An ideal measurement is a matching rule that yields higher matching score for the image pairs belonging to the same person than the pairs belonging to different persons, which can be a big challenge if images are captured by different views/cameras with different setups and illumination conditions (*i.e.*, a typical scenario found in person re-identification, usually not handled by direct distance metric comparison). As highlighted in [5], similarity measurements which are learned (*e.g.*, [6, 7]) from training samples generally enjoy better accuracy performance than learning free methods [8]. Note that the goal of metric learning algorithms is to take advantage of prior information in form of labels over simpler though more general similarity measures [9]. The achieved results are then provided in the form of a list of ranked matching persons. It often happens that the true match is not ranked first but it is in the first positions. This is mostly due to the visual ambiguities shared between the true match and other "similar" persons [10].

In order to address the re-identification problem, existing methods exploit either feature representation [11, 12, 13] or metric learning [9, 7]. In feature representation, robust and discriminative features are constructed

---
[*]Corresponding author
*Email address:* juliojj@gmail.com (Julio C. S. Jacques Junior)

such that they can be used to describe the appearance of the same individual across different camera views under various conditions [14], whereas distance metric learning methods attempt to learn a metric in the space defined by image features that keep features coming from same class closer, while, the features from different classes are farther apart [2]. Recently, Convolutional Neural Networks (CNN) have been adopted in person re-identification [15, 11], providing a powerful and adaptive tool to handle computer vision problems without excessive usage of hand-crafted image features. However, as mentioned in the work of Wu et al. [11], hand-crafted concatenation of different appearance features sometimes would be more distinctive and reliable, due to significant changes in view angle, lighting, background clutter and occlusion.

In this work we exploit the best of different state-of-the-art models to advance the field of person re-identification. The proposed model is inspired by the work of Chen et al. [5], which enforces similarity learning with spatial constraints, and achieved (by the time of its publication) the best score (i.e., top rank recognition rate) on VIPeR [16] dataset (which is one of the most challenging datasets employed in person re-identification). In this paper, by combining new and complementary features within [5], followed by a post-ranking [10] and a ranking aggregation strategy [17], we advance the state-of-the-art in person re-identification on two public datasets, VIPeR and PRID450s [18] (by 2.43% and 2.66%, respectively) as well as achieve competitive results on CUHK01 [19] dataset.

The new and complementary adopted features can be briefly enumerated as follows: (i) Salient Color Names based Color Descriptor (SCNCD) [6] combined with color histogram (to encode color information), Histogram of Oriented Gradients (HOG) [20] and Scale Invariant Local Ternary Patterns (SILTP) [21] (to encode texture information). Although HOG and SILTP were exploited in [5], they were not combined with SCNCD; (ii) SCNCD combined with background/foreground information, automatically extracted via Deep Decompositional Network (DDN) [22]; (iii) Gaussian Of Gaussian (GOG) descriptor [23], which encodes both color and texture information; (iv) Convolutional Neural Network (CNN) features constrained by hand-crafted color histograms [11] and combined with Local Maximal Occurrence (LOMO) features [24]. A quantitative analysis regarding the effectiveness of each complementary feature is presented on Sec. 4.4. Experimental results showed that the proposed new features demonstrated to complement each other, being very powerful when combined with a ranking aggregation strategy.

The rest of the paper is organized as follows: Section 2 presents the state-of-the-art concerning person re-identification. The proposed model is described in Section 3, and experimental results are provided in Section 4. Finally, conclusions are given in Section 5.

## 2. RELATED WORK

Existing research on person re-identification has concentrated either on the development on sophisticated and robust features to describe the visual appearance of a person under significant visual variabilities or on the development of new learning distance metrics. In this section we present the state-of-the-art on person re-identification, briefly describing the works that achieved the best recognition rates on three broadly employed public datasets, VIPeR, PRID450s and CUHK01, without focusing on the standard taxonomy (i.e., feature representation or metric learning).

As in the work of Paisitkriangkrai et al. [14], one simple approach to exploit multiple visual features is to build an ensemble of distance functions, in which each distance function is learned using a single feature and the final distance is calculated from a weighted sum of these distance functions. However, the usage of predetermined weights is undesirable as highly discriminative features in one environment might become irrelevant in another one. In their work, a model to learn weights of these distance functions by optimizing the relative distance or by maximizing the average rank-k recognition rate is proposed. Mirmahboub et al. [25] proposed a novel re-ranking method based on a fusion scheme that reweights an ensemble of distance metric outcomes according to their discriminative capacity. They particularly show that the fused distance perform largely better than any of the distances inferred by each feature separately.

To consider spatial information, a common usage in person re-identification is to divide the person image into few regions/stripes and concatenate dense local features to implicitly encode the spatial layout of the person. Chen et al. [5] proposed a model for person re-identification that combines spatial constraints and the polynomial feature map [7] into a unified framework. They mention that enforcing the matching within corresponding regions can effectively reduce the risk of mismatching and become more robust to partial occlusions. In addition, their framework can benefit from the complementarity of global and local similarities.

The post-ranking method for person re-identification is a relatively unexplored area [10] which has been attracting a lot of attention from the research community. Prates and Schwartz [17] presented a Color-based Ranking Aggregation (CBRA) method, which explores different feature representations to obtain complementary ranking lists, and combine them in order to improve person re-identification. In their work, the KISSME [9] metric learning was adopted and different strategies for ranking aggregation, based on the Stuart rank aggregation method [26], were proposed. García et al. [27, 10] related that inspections on the ranked matches can be applied to refine the output in such a way that the correct match will have higher probability to be found in the first ranks. Hence, their work is founded on the idea that a ranking, achieved by any algorithm, con-

tains valuable information which can be further exploited to improve the rank of the true match. To achieve such a goal, they propose an unsupervised post-ranking framework. Once the initial ranking is available, content and context sets are extracted. Then, these are exploited to remove the visual ambiguities and to obtain discriminant feature space which is finally exploited to compute the new ranking.

Bai et al. [28] studied person re-identification with manifold-based affinity learning. In their work, a novel affinity learning algorithm called Supervised Smoothed Manifold (SSM) is proposed, which can be plunged into most existing algorithms, serving as a generic postprocessing procedure to further boost identification accuracy.

In relation to domain adaptation in machine learning, Chen et al. [12] proposed a schema called Mirror Representation to address the view-specific feature distortion problem in person re-identification. It embeds the view-specific feature transformation and enables alignment of the feature distributions across disjoint views for the same person. Zhang and collaborators [29] argue that most existing approaches focus on learning a fixed distance metric for all instance pairs, while ignoring the individuality of each person. They formulate person re-identification as an imbalanced classification problem and learn a classifier specifically for each pedestrian such that the matching model is highly tuned to the individual appearance.

Considering the recently proposed CNN based methods for person re-identification, in [11] a deep Feature Fusion Network (FFN) is proposed in order to use hand-crafted features to regularize CNN process so as to make the convolutional neural network extract features complementary to hand-crafted ones. As mentioned by the authors, different to other deep methods for person re-identification (*e.g.*, [15, 30]) which are based on pairwise input, they can directly extract deep features on single images, being able to be learnt by any conventional classifier. Xiao et al. [13] presented a pipeline for learning deep feature representations from multiple domains with CNN. Authors argue that when training a CNN with data from all domains, some neurons learn representations shared across several domains, while some others are effective only for a specific one. Based on this observation they proposed a Domain Guided Dropout algorithm (a method of muting non-related neurons for each domain). Liu et al. [31] proposed a new soft attention-based model, *i.e.*, the end-to-end Comparative Attention Network (CAN), specifically tailored for the task of person re-identification, which can adaptively find multiple local regions with discriminative information in person images in a recurrent way. Such approach learns to selectively focus on parts of pairs of person images after taking a few glimpses of them and adaptively *comparing* their appearances.

Although a large number of existing approaches have exploited state-of-the-art visual features, advanced metric learning algorithms, post-ranking or ranking aggregation strategies, domain adaptation based models or even CNN based ones, state-of-the-art results on commonly evaluated person re-identification benchmarks is still far from the accuracy performance needed for most real-world surveillance applications [14].

## 3. PROPOSED MODEL

In this work, we propose to exploit different feature representations[1] to advance the state-of-the-art in person re-identification. In the proposed model, each image is represented in different ways, which include hand-crafted descriptors and deep features. To describe the matching between a probe image and a gallery set, a similarity learning metric built on the polynomial feature map [7] is adopted, also taking into account spatial (local and global) information. As each image has different descriptors, different similarities are computed, according to each representation. This way, for each probe image and gallery set, different rank lists are generated, each one assigned to each feature representation. Once these initial rankings are available, content and context information[10] are extracted for each feature representation and respective probe image. Then, these are exploited to remove the visual ambiguities and to obtain discriminant feature space which is finally exploited to compute new ranking lists. The final rank list is obtained through ranking aggregation, which combines these complementary ranking lists. An overview of our model is illustrated in Fig. 1.

Next, we briefly revisit the polynomial feature map and the spatially constrained techniques [5][2], as they are the basis of the proposed model. In a second stage, we describe the proposed complementary features. Finally, we describe the adopted post-ranking and ranking aggregation strategies.

### 3.1. Polynomial Feature Map

In order to measure the similarity between image descriptors $\mathbf{x}_a, \mathbf{x}_b \in \mathbb{R}^{d \times 1}$, we learn the similarity function as:

$$f(\mathbf{x}_a, \mathbf{x}_b) = \langle \phi(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W} \rangle_F, \qquad (1)$$

where $\langle \cdot, \cdot \rangle_F$ is the Frobenius inner product. To take advantage of both Mahalanobis distance and bilinear similarity metric, we decompose $f(\mathbf{x}_a, \mathbf{x}_b)$ as follows:

$$f(\mathbf{x}_a, \mathbf{x}_b) = \langle \phi_M(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}_M \rangle_F + \langle \phi_B(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}_B \rangle_F. \qquad (2)$$

The part $\langle \phi_M(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}_M \rangle_F = (\mathbf{x}_a - \mathbf{x}_b)^\top \mathbf{W}_M (\mathbf{x}_a - \mathbf{x}_b)$ is connected to the Mahalanobis distance. The part $\langle \phi_B(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}_B \rangle_F = \mathbf{x}_a^\top \mathbf{W}_B \mathbf{x}_b + \mathbf{x}_b^\top \mathbf{W}_B \mathbf{x}_a$ corresponds to bilinear similarity. Both parts ensure the effectiveness

---

[1] An evaluation about different color spaces and their combinations for person re-identification can be found in [32].
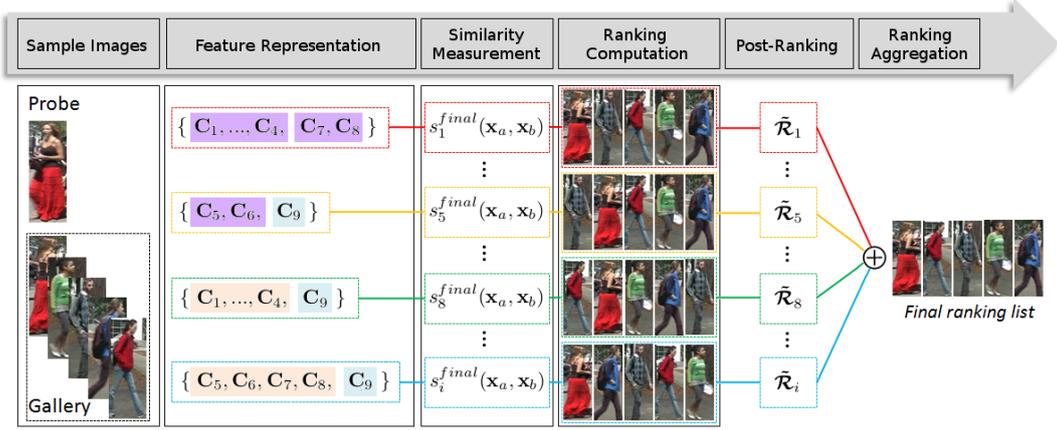[2] Code available at https://github.com/dapengchen123/SCSP

Figure 1: Overview of the proposed model. For each sample image, different visual cues are defined (*i.e.*, $\{\mathbf{C}_1, ..., \mathbf{C}_9\}$, as detailed in Sec. 3.3). Then, different feature representations are proposed, taking into account global (cyan regions), local (salmon regions) or both global and local (pink regions) information. For each probe image and gallery set, different similarity measures are computed using different feature representations. Each representation produces a initial ranking list based on the adopted similarity function. Then, a post-ranking approach is used in order to improve the recognition rate obtained by each initial ranking list. The final ranking list is obtained through ranking aggregation, which combines complementary ranking lists obtained from different feature representations.

of $f(\mathbf{x}_a, \mathbf{x}_b)$. The dimensionality of the feature map is reduced by means PCA for $\mathbf{x}_a$ and $\mathbf{x}_b$ before its generation[3].

### 3.2. Spatially Constrained Similarity Function

#### 3.2.1. Regional feature map

The input image is partitioned into $R$ non-overlap horizontal stripe regions. Each region is divided into a collection of overlapped patches, from which we extract color and texture histograms. The extracted histograms belonging to a same stripe region are concatenated together. After that, PCA is applied to reduce the dimensionality and to obtain the region descriptor $\mathbf{x}^r$ for the $r$-th stripe, where $r \in \{1, ..., R\}$. A stripe region $r$ can be described by $C$ visual cues $\{\mathbf{x}^{r,1}, ..., \mathbf{x}^{r,c}, ..., \mathbf{x}^{r,C}\}$, thus $\mathbf{x}_a$ and $\mathbf{x}_b$ accordingly form $C$ polynomial feature maps for the $r$-th region, *i.e.*, $\{\phi^{r,1}(\mathbf{x}_a, \mathbf{x}_b), ..., \phi^{r,c}(\mathbf{x}_a, \mathbf{x}_b), ..., \phi^{r,C}(\mathbf{x}_a, \mathbf{x}_b)\}$, where $\phi^{r,c}(\mathbf{x}_a, \mathbf{x}_b) = \phi(\mathbf{x}_a^{r,c}, \mathbf{x}_b^{r,c})$.

#### 3.2.2. Local similarity integration

In order to exploit the complementary strengths of multiple visual cues within a local region, a linear similarity function is employed to combine them together for the $r$-th region:

$$s^r(\mathbf{x}_a, \mathbf{x}_b) = \sum_{c=1}^{C} \langle \phi^{r,c}(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}^{r,c} \rangle_F, \qquad (3)$$

where $\mathbf{W}^{r,c} = [\mathbf{W}_M^{r,c}, \mathbf{W}_B^{r,c}]$ and $\mathbf{W}_M^{r,c}$, $\mathbf{W}_B^{r,c}$ correspond to $\phi_M^{r,c}(\mathbf{x}_a, \mathbf{x}_b)$ and $\phi_B^{r,c}(\mathbf{x}_a, \mathbf{x}_b)$, respectively. The local similarities scores are integrated as:

$$s^{local}(\mathbf{x}_a, \mathbf{x}_b) = \sum_{r=1}^{R} s^r(\mathbf{x}_a, \mathbf{x}_b). \qquad (4)$$

#### 3.2.3. Global-local collaboration

In order to describe the matching of large patterns across the stripes, the polynomial feature map is also used for the whole image, yielding global similarity:

$$s^{global}(\mathbf{x}_a, \mathbf{x}_b) = \sum_{c=1}^{C} \langle \phi^{G,c}(\mathbf{x}_a, \mathbf{x}_b), \mathbf{W}^{G,c} \rangle_F, \qquad (5)$$

where $\mathbf{W}^{G,c} = [\mathbf{W}_M^{G,c}, \mathbf{W}_B^{G,c}]$ and $\mathbf{W}_M^{G,c}$, $\mathbf{W}_B^{G,c}$ correspond to $\phi_M^{G,c}(\mathbf{x}_a, \mathbf{x}_b)$ and $\phi_B^{G,c}(\mathbf{x}_a, \mathbf{x}_b)$, respectively. Here, $\phi^{G,c}(\mathbf{x}_a, \mathbf{x}_b) = \phi(\mathbf{x}_a^{G,c}, \mathbf{x}_b^{G,c})$ and $\mathbf{x}_a^{G,c}, \mathbf{x}_b^{G,c}$ are the $c$-th type global visual descriptors for image $a$ and $b$. Finally, the global similarity and local similarity are linearly combined, and the overall similarity score is given by:

$$s^{final}(\mathbf{x}_a, \mathbf{x}_b) = s^{local}(\mathbf{x}_a, \mathbf{x}_b) + \gamma s^{global}(\mathbf{x}_a, \mathbf{x}_b), \quad (6)$$

where $\gamma$ is the hyper-parameter that mediates the local and global similarities (experimentally set to $\gamma = 1.1$).

#### 3.2.4. Visual Cues and Parameter settings

In the original model of [5], four visual cues are used (*i.e.*, $C = 4$). First, images are resized to 48×128. Each region $r$ (from $R$, experimentally set to $R = 4$)[4] is divided into a set of local patches (with 8×16 of size and stride of 4×8). For each patch, six types of features are extracted: $HSV_1$, $LAB_1$ (are 8×8×8 joint histograms), $HSV_2$, $LAB_2$ (are 48 bin concatenated histograms with each channel having 16 bins), HOG [20] and SILTP [21] (texture descriptors). The four visual cues $\mathbf{C}_1$, $\mathbf{C}_2$, $\mathbf{C}_3$ and $\mathbf{C}_4$ concatenate both color and texture features, which are organized as $HSV_1$/HOG, $HSV_2$/SILTP, $LAB_1$/SILTP and $LAB_2$/HOG, respectively.

---

[3]A detailed explanation about how $\mathbf{W}_M$ and $\mathbf{W}_B$ are learned using the ADMM optimization algorithm can be found in [5].

[4]A default $R$ value was adopted from [5].

Regarding each visual cue, descriptors generated for each patch, within a specific region $r$, are concatenated to compose the descriptor of such region. Similarly, the global descriptor is generated through the concatenation of the descriptors computed for all patches. For each visual cue, obtained color and texture descriptors are normalized (to have unit $L_2$ norm) before concatenation. Then, each visual cue is reduced by PCA. Finally, the resulting descriptor is normalized again in the same way. As mentioned in [5], the PCA reduced dimension $d$ depends on the size of training data. In our experiments we adopted $d$ to be 120 for all evaluated datasets.

### 3.3. Complementary features

In order to improve state-of-the-art recognition performance in person re-identification, we propose to include new and complementary features within the similarity function presented in [5], as described next.

### 3.3.1. SCNCD [6]

For each color to be named, salient color names indicate that a color only has a certain probability of being assigned to several nearest color names, and that the closer the color name is to the color, the higher probability the color has of being assigned to this color name. Through this way, we can assign multiple similar colors to the same index with the same color descriptor.

Color distributions over color names in different color spaces are then obtained and fused to generate a feature representation. In addition, and similarly to [6], color histogram is computed for each color channel and fused with color names distribution (the number of bins is set to 32). In this work, SCNCD and color histograms are extracted using the original RGB, normalized $rgb$, $l_1l_2l_3$ and HSV color models. Such procedure is performed locally, regarding each region $r$, as well as globally, regarding the whole image. To be specific, SCNCD are extracted similarly to [6], except that in our model the image is divided in 4 regions ($R = 4$) and a second subdivision at the local level is performed. First, the image is divided into $R$ horizontal stripes, from which features are extracted and concatenated (global descriptor), as illustrated on the right side of Figure 2. Then, each region $r$ is subdivided again into $R$ horizontal stripes, from which features are extracted and concatenated (local descriptor), as illustrated on the left side of Figure 2.

Two new visual cues are then proposed, $\mathbf{C}_5$ and $\mathbf{C}_6$. Both concatenate color and texture features, which are organized as SCNCD/HOG and SCNCD/SILTP, respectively. In this case, HOG and SILTP are extracted in the same way as in [5]. As before, obtained descriptors that compose each new visual cue are normalized (to have unit $L_2$ norm) before final concatenation. The resulting descriptor is then reduced by means of PCA before final normalization step.
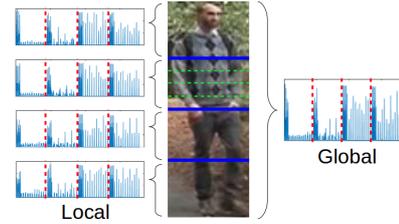


Figure 2: Illustration of the generated global and local descriptors based on SCNCD (and color histograms).

### 3.3.2. Background/foreground information

Due to the fact that the background in person re-identification is not constant and may even include disturbing factors, background feature representation combined directly with the foreground feature representation may reduce classification accuracy. To address this problem, [6] proposed an image-foreground feature representation, which can be seen as that the foreground information is employed as the main information while the background information is treated as a secondary one. Differently from [6], we propose to extract the foreground mask with a more powerful segmentation model based on Deep Decompositional Network (DDN) [22].

The DDN was developed to tackle the problem of pedestrian parsing, and designed to segment pedestrian images into semantic regions, such as hair, head, body, arms, and legs. It directly maps low-level visual features (HOG) to the label maps of body parts, being able to accurately estimate complex pose variations while being robust to occlusions and background clutters. In a nutshell, DDN jointly estimates occluded regions and segments body parts by stacking three types of hidden layers: occlusion estimation layers, completion layers, and decomposition layers. The occlusion estimation layers estimate a binary mask, indicating which part of a pedestrian is invisible. The completion layers synthesize low-level features of the invisible part from the original features and the occlusion mask. The decomposition layers directly transform the synthesized visual features to label maps. Fig. 3 illustrates some binary masks automatically obtained using [22][5].



Figure 3: Input images and respective binary masks obtained using [22].

---

[5]Code available at `http://mmlab.ie.cuhk.edu.hk/projects/luoWTiccv2013DDN/`

### 3.3.3. Gaussian Of Gaussian (GOG)

Matsukawa et al. [23] proposed a region descriptor based on hierarchical Gaussian distribution of pixel feature. In their work, local patches inside a region are densely extracted and the region is regarded as a set of local patches. The region is modeled as a set of multiple Gaussian distributions, each of them representing the appearance of one local patch. The characteristics of the set of patch Gaussians are again described by another Gaussian distribution (defined as *a region Gaussian*). The parameters of the region Gaussian are then used as feature vector to represent the region. The GOG descriptor provides a consistent way to generate discriminative and robust features that describe color and textural (*e.g.*, gradient magnitudes along different directions) information simultaneously.

In our work, we adopted the $GOG_F$ (Fusion) descriptor[6], extracted as described in [23], which concatenates different GOG descriptors generated from different colorspaces (RGB, LAB, HSV and normalized *rgb*). Thus, a new visual cue is proposed, $\mathbf{C}_7$. As before, obtained descriptor is normalized (to have unit $L_2$ norm). The resulting descriptor is then reduced by means of PCA before final normalization. Note that, as the final representation is a concatenation of local features, it will be used in Eq. 6 as a global descriptor (similarly as $\mathbf{C}_8$, described next).

### 3.3.4. Deep feature [11]

Feature Fusion Net (FFN) is used to allow deep feature representation in the adopted framework, as it demonstrated to be very effective in person re-identification tasks. FFN consists of two parts. The first part deals with traditional convolution, pooling and activation neurons for input images. The second part of the network processes additional hand-crafted feature representations of the same image. Both, CNN features and the hand-crafted features are followed by a fully connected layer and then linked together in order to produce a full-fledge image description from the last convolutional layer.

Regarding the hand-crafted features, authors first modified the Ensemble of Local Features (ELF) [33] by improving the color space and stripe division (denoted as ELF16). Input images are equally partitioned into 16 horizontal stripes, and the features are composed of color features including RGB, HSV, LAB, XYZ, YCbCr and NTSC, and texture features including Gabor, Schimid and LBP. A 16D histogram is extracted for each channel and then normalized by $L_1$ norm. All histograms are concatenated together to form a single vector. The FFN was then trained on the Market-1501 [34] dataset, which is the largest public person re-identification dataset up to date, composed of 38195 images from 1501 identities.

The authors of [11] also mention that even though the proposed CNN-based feature performs better when com-

pared to LOMO [24] features, the combination of both kind of features demonstrates to have higher discriminative power. Thus, the concatenation of both (CNN-based feat.+LOMO) is defined in their work as the final representation (denoted in our work, from now, by just *Deep feature*. We also apply PCA (as previously mentioned), to reduce the dimensionality of the resulting Deep feature, which is then normalized by $L_2$ norm. This final representation is used as another complementary cue ($\mathbf{C}_8$). Note that $\mathbf{C}_8$ composes a representation for the whole image, so it will be only used as a global descriptor.

### 3.3.5. Integrating complementary features

To integrate the new and complementary features, we compute different similarity measures using different feature representations (for each pair of images being compared), which are then exploited next by the post-ranking and ranking aggregation strategies. To be specific, we compute $s_i^{final}(\mathbf{x}_a^{F_i}, \mathbf{x}_b^{F_i})$, defined in Eq. 6, where $F_i$ are different feature representations, described in Table 1. We refer $\{F_7, ..., F_{12}\}$ as simplified versions of $\{F_1, ..., F_6\}$.

Table 1: Summary of the adopted feature representations $F_i$ ($F_0$ is the baseline [5]). G, L and GL indicate, respectively, if the visual cue is applied just on the global, local or both parts of the Eq. 6.

|  | baseline cues | | | | scncd | | gog | deep |
|---|---|---|---|---|---|---|---|---|
|  | $\mathbf{C}_1$ | $\mathbf{C}_2$ | $\mathbf{C}_3$ | $\mathbf{C}_4$ | $\mathbf{C}_5$ | $\mathbf{C}_6$ | $\mathbf{C}_7$ | $\mathbf{C}_8$ |
| $F_0$ | GL | GL | GL | GL | - | - | - | - |
| $F_1$ | GL | GL | GL | GL | - | - | G | - |
| $F_2$ | GL | GL | GL | GL | - | - | - | G |
| $F_3$ | GL | GL | GL | GL | - | - | G | G |
| $F_4$ | - | - | - | - | GL | GL | G | - |
| $F_5$ | - | - | - | - | GL | GL | - | G |
| $F_6$ | - | - | - | - | GL | GL | G | G |
| $F_7$ | L | L | L | L | - | - | G | - |
| $F_8$ | L | L | L | L | - | - | - | G |
| $F_9$ | L | L | L | L | - | - | G | G |
| $F_{10}$ | - | - | - | - | L | L | G | - |
| $F_{11}$ | - | - | - | - | L | L | - | G |
| $F_{12}$ | - | - | - | - | L | L | G | G |

### 3.4. Post-ranking based on DCIA

According to Garcia et al.[27, 10], additional ranking inspections on the ranked matches can be applied to refine the output in such a way that the correct match will have higher probability to be found in the first ranks. To this end, they proposed the Discriminant context information Analysis (DCIA) method, which is built under the definition of content and context information. The content information is the set of gallery images that have low dissimilarity with respect to the probe. The context information is the set of gallery images that have low dissimilarity with both the probe and an image of the content information. In this subsection we introduce basic concepts related to their method as well as describe how post-ranking is performed.

---

[6]Available at http://www.i.kyushu-u.ac.jp/~matsukawa/ReID.html

### 3.4.1. Definitions

Let $\mathcal{A} = \{\mathbf{I}_p^A\}_{p=1}^N$ be the set of $N$ probe images and $\mathcal{B} = \{\mathbf{I}_g^B\}_{g=1}^M$ be the set of $M$ gallery images. Given a probe image $\mathbf{I}_p^A$, its initial ranking is defined as $\mathcal{R}_p = \{\mathbf{I}_i^B\}_{i=1}^M$, where the gallery images $\mathbf{I}_i^B$ are sorted depending on the dissimilarity to the probe. In other words, $d(\mathbf{I}_p^A, \mathbf{I}_i^B) < d(\mathbf{I}_p^A, \mathbf{I}_{i+1}^B)$, where $d(\cdot, \cdot)$ is a suitable dissimilarity measure (i.e., as defined in Eq. 6) and $i$ goes from 1 to $M - 1$. $\boldsymbol{\mathcal{R}} = \{\mathcal{R}_p^N\}_{p=1}^N$ denotes the set of such initial rankings computed for the $N$ probes.

### 3.4.2. Content Information

The content information is defined as the set of features extracted from the correlated matches, i.e., a subset of gallery images $\mathcal{B}^{cn} \subseteq \mathcal{B}$ present in the fist ranks and which are likely to share visual ambiguities with the probe. Elements in such a set are selected from the top $m$ positions in the initial ranking $\mathcal{R}_p$ which have matching distance less than a specific threshold $Th$. The $m$ value, as well as the adopted threshold, are dynamically computed for each probe image, based on the shape of dissimilarities vs rank plots (see [27] for additional details). Thus, the set of $m$ correlated matches equals $\mathcal{B}^{cn} = \{\mathbf{I}_i^B | d(\mathbf{I}_p^A, \mathbf{I}_i^B) \leq Th\}$. Therefore, the content set $\mathcal{C}_p^{cn} = \{\mathbf{x}_1^{cn}, ..., \mathbf{x}_m^{cn}\}$ contains the $m$ feature vectors extracted from the correlated matches in $\mathcal{B}^{cn}$. Notice that, only images in $\mathcal{C}_p^{cn}$ are re-ranked.

### 3.4.3. Context Information

The context information is given by the $K$-common nearest neighbors of the probe and a correlated match. Given $\mathbf{I}_p^A$, its respective context set is extracted by exploiting $\mathcal{C}_p^{cn}$. First, the initial rank list $\mathcal{R}_g$ is computed for each correlated matching image $\mathbf{I}_g^B \in \mathcal{C}_p^{cn}$ by evaluating its similarity with images in the gallery set $\mathcal{B}^* = (\mathcal{B} \backslash \mathbf{I}_g^B)$ using model parameters ($\phi$ and $\mathbf{W}$) and distance $d(\cdot, \cdot)$, i.e., as defined in Eq. 6. Then, given $\mathcal{R}_g$, we compute the top $m_g$ positions which have matching distance less than $Th_g$ (being $m_g$ and $Th_g$ computed in the same way as $m$ and $Th$, respectively). These elements represent the images that have high similarity with both the probe $\mathbf{I}_p^A$ and the correlated match $\mathbf{I}_g^B$. The context information is extracted from the $K$-common context matches. Feature vectors extracted from such images form the context information set $\mathcal{C}_p^{cx} = \{\mathbf{x}_1^{cx}, ..., \mathbf{x}_n^{cx}\}$, where $n = K$. Finally, $\mathcal{C}_p^{cx}$ is updated by removing images that are in duplicate with $\mathcal{C}_p^{cn}$. The hard threshold $K$ was set experimentally to $K = 13$ as in [27].

Nevertheless, we observed the respective $K$-common context matches are obtained, for some probe images, from a flat histogram. In this case, the $K$-common context matches might be imprecisely obtained, mainly when the number of images that compose the histogram are greater than $K$. Thus, we introduced a new condition to also consider the similarity from the correlated context match

to the probe image in order get the $K$-common context matches most similar to $\mathbf{I}_p^A$.

### 3.4.4. Discriminative Information Analysis

Given a probe image $\mathbf{I}_p^A$, let $\mathcal{D}_p = \{\mathbf{x}_p, \mathcal{C}_p^{cn}, \mathcal{C}_p^{cx}\}$ be the set composed of its feature vector and of feature vectors obtained in the content and context information. $\mathcal{D}_p$ is redefined as a feature matrix $\mathbf{D}_p \in \mathbb{R}^{d \times l}$ with zero mean, where $l = 1 + m + n$ is the number of vectors. Let $\mathbf{P} \in \mathbb{R}^{d \times k}$ be the first $k$ components of $\mathbf{D}_p$ selected to represent the common appearance subspace. Thus, the discriminant information can be obtained as

$$\mathbf{D}_p^* = \mathbf{D}_p - \mathbf{P}\mathbf{P}^T \mathbf{D}_p, \tag{7}$$

where each column of $\mathbf{D}_p^*$ represents a discriminant feature vector $\mathbf{x}^*$. Differently from [27], where $k$ principal components corresponding to the $55\%$ ($k = 0.55$) of energy of the set of feature vectors have been used to represent the common appearance subspace, in this work we empirically defined $k = 0.35$.

### 3.4.5. Re-ranking Training

The DCIA is first applied to the train set $\mathcal{I}_{Tr}$. More specifically, it is applied to each ranking $\mathcal{R}_p^{Tr} \in \boldsymbol{\mathcal{R}}^{Tr}$. As result, the discriminant feature vectors $\mathbf{x}_p^{*A}$ and $\mathbf{x}_g^{*B} \in \mathbf{D}_p^{*Tr}$ are obtained for each probe image $p$. The resulting sets $\mathbf{x}_{Tr}^{*A}$ and $\mathbf{x}_{Tr}^{*B}$ together with the pairwise labels are used to learn the new model parameters $\phi^*$ and $\mathbf{W}^*$.

### 3.4.6. Post-ranking Optimization

Given a test rank in $\boldsymbol{\mathcal{R}}$, the DCIA is performed to obtain the discriminative test feature vectors $\mathbf{x}_p^{*A}$ and $\mathbf{x}_g^{*B}$. Then, the set of such vectors $\{\mathbf{x}^{*A}, \mathbf{x}^{*B}\}$ is evaluated by the new model parameters $\phi^*$ and $\mathbf{W}^*$. The obtained distances are used to re-rank the correlated matches, hence to compute the final ranking $\tilde{\boldsymbol{\mathcal{R}}}$. Such procedure is performed for each feature representation.

### 3.5. Ranking Aggregation Strategy

We propose to explore different feature representations to obtain complementary ranking lists and combine them using the Stuart ranking aggregation method [26]. The Stuart ranking aggregation method, which was originally designed to define a gene-coexpression network over DNA microarrays from humans, flies, worms, and yeast, is a probabilistic method based on order statistics to evaluate the probability of observing a particular configuration of ranks across the different organisms, even when there are irrelevant and noise inputs. The significance of the interactions in the network is verified by means of a variety of statistical tests [7].

---

[7] An optimized solution of [26] is presented in [35].

Let first denote by $\oplus$ the aggregation operator, for instance if $\tilde{\mathcal{R}}_n = \tilde{\mathcal{R}}_1 \oplus \tilde{\mathcal{R}}_2 \oplus ... \oplus \tilde{\mathcal{R}}_{n-1}$, then $\tilde{\mathcal{R}}_n$ is a ranking list computed by the aggregation of ranking lists from $\tilde{\mathcal{R}}_1$ to $\tilde{\mathcal{R}}_{n-1}$. As we use different descriptors to represent each image, and have adopted a strategy in which we can measure the similarity $s_i^{final}(\mathbf{x}_a, \mathbf{x}_b)$ of image pairs using different ways (Sec. 3.3.5), we are also able to compute different ranking lists for each probe image and gallery set, as illustrated in Fig. 1. Moreover, a *tunning* strategy is performed on the final aggregation in order to improve the results, as explained next.

Consider the list of complementary features $\mathcal{L} = \{F_1, ..., F_{12}\}$, sorted based on average top-1 rank recognition rate (obtained from a validation set), where $\mathcal{L}_i$ has better accuracy than $\mathcal{L}_{i+1}$. Thus, we can aggregate the ranking lists that obtained higher accuracies (*i.e.*, the best-$n$ feature representations) and ignore "poor" ranking lists that may push down the final result/aggregation. Such procedure is performed by the aggregation of the ranking lists of the best-2 feature representations (*i.e.*, $\{\mathcal{L}_1, \mathcal{L}_2\}$), best-3 ($\{\mathcal{L}_1, ..., \mathcal{L}_3\}$), and so on, up to the aggregation of the whole list (*i.e.*, best-12, $\{\mathcal{L}_1, ..., \mathcal{L}_{12}\}$).

# 4. EXPERIMENTAL RESULTS

In order to demonstrate the effectiveness of the proposed model, this section presents experimental results on three broadly employed public datasets for person re-identification, *i.e.*, VIPeR [16], PRID450s [18] and CUHK01 [19]. Five case studies were performed. First, (i) the proposed model was compared against state-of-the-art person re-identification models using a well known evaluation protocol (Sec. 4.2). Then, we decomposed the proposed complementary features and performed the following experiments: (ii) influence of the background/foreground information within SCNCD (Sec. 4.3); (iii) accuracy performance obtained by each complementary feature (Sec. 4.4); (iv) improvements obtained by post-ranking (Sec. 4.5). Finally, (v) the best-$n$ *tunning* strategy for rank-aggregation was analyzed (Sec. 4.6).

The adopted datasets are presented in two disjoint camera views, with significant misalignment, light changes and body part distortion. Table 2 summarizes the three datasets. Challenging image samples (due to illumination problems, pose variation, occlusions or even by high similarity between different people) are illustrated in Fig. 4.

Table 2: Summary of the adopted datasets.

|  | VIPeR | PRID450s | CUHK01 |
|---|---|---|---|
| Images | 1264 | 900 | 3884 |
| Individuals (ID) | 632 | 450 | 971 |
| Images per ID (per view) | 1 | 1 | 2 |



| (a) VIPeR | (b) PRID450s | (c) CUHK01 |

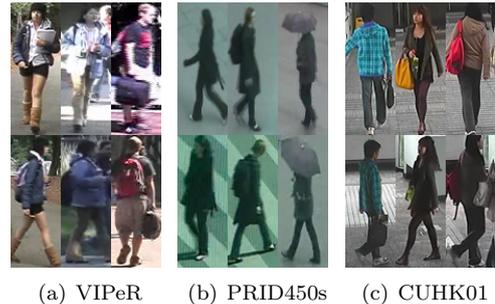Figure 4: Sample images of the adopted datasets. Images on the same column represent the same person.

## 4.1. Evaluation Protocol

Our experiments follow the evaluation protocol defined in [14] for a single-shot scenario, *i.e.* we randomly partitioned each dataset into two parts, 50% for training and 50% for testing, without overlap on person identities. As the CUHK01 dataset contains 971 individuals, 485 of them were randomly sampled for training and the rest for testing, as in [29]. Images from camera A are used as probe and those from camera B as gallery. For the CUHK01 dataset, in which each individual has two images per camera view, we randomly selected one image of the individual taken from the camera A as the probe image and one image of the same individual taken from the camera B as the gallery image. For all evaluated datasets, each probe image is matched with every image in gallery and the rank of correct match is obtained. This procedure is repeated 10 times and the average of Cumulative Matching Characteristic (CMC) curves across 10 partitions is reported.

## 4.2. Case 1: State-of-the-art comparison

This experiment compares the overall accuracy performance of the proposed model in relation to the state-of-the-art. Different feature representations were integrated, as described in Sec. 3.3.5, followed by the post-ranking approach (described in Sec.3.4) and ranking aggregation strategies (described in Sec. 3.5). Table 3 summarizes the obtained results.

As it can be seen in Table 3, the proposed model outperforms the state-of-the-art on both VIPeR and PRID450s datasets, and achieved competitive results on CUHK01 dataset. Some other works obtained better results than ours on CUHK01 dataset, however, they were not included in this comparison as they either use a different evaluation protocol [36] or include additional data in the train set [37] (*i.e.*, CUHK03 database, which was captured in the same environment as CUHK01 and could benefit when CUHK01 is evaluated as both share similar features).

We can also observe the CAN-VGG16 method [31], which obtained promising results on CUHK01 dataset, were outperformed by the proposed model on VIPeR dataset by a significant margin. The slow performance

Table 3: State-of-the-art comparison. Top Matching Rank (%) on the three adopted datasets.

| Rank | 1 | 5 | 10 | 20 |
|------|-----|-----|-----|-----|
| **VIPeR** | | | | |
| **Our best**-10 | **67.21** | **87.78** | 93.39 | **97.82** |
| **Our best**-12 | 66.83 | 87.78 | **93.41** | 97.72 |
| DCIA [10] | 64.78 | 76.85 | 86.88 | 94.77 |
| Re-ranking [25] | 59.46 | 86.68 | 93.39 | 97.63 |
| SSM [28] | 53.73 | - | 91.49 | 96.08 |
| SCSP [5] | 53.54 | 82.59 | 91.49 | 96.65 |
| Deep+LOMO [11] | 51.06 | 81.01 | 91.39 | 96.90 |
| CAN-VGG16 [31] | 47.20 | 79.20 | 89.20 | 95.80 |
| Mirror [12] | 42.97 | 75.82 | 87.28 | 94.84 |
| LSSCDL [29] | 42.66 | - | 84.27 | 91.93 |
| **PRID450s** | | | | |
| **Our best**-3 | **75.64** | **93.38** | 96.44 | 98.22 |
| **Our best**-12 | 73.91 | 92.58 | 95.87 | 97.87 |
| SSM [28] | 72.98 | - | **96.76** | **99.11** |
| Deep+LOMO [11] | 66.62 | 86.84 | 92.84 | 96.89 |
| LSSCDL [29] | 60.49 | - | 88.58 | 93.60 |
| Mirror [12] | 55.42 | 79.29 | 87.82 | 93.87 |
| **CUHK01** | | | | |
| CAN-VGG16 [31] | **67.20** | 87.30 | **92.50** | **97.20** |
| **Our best**-3 | 66.91 | 86.95 | 92.12 | 95.7 |
| LSSCDL [29] | 65.97 | $\approx$ **88.0** | $\approx$ 92.0 | $\approx$ 96.0 |
| **Our best**-12 | 64.28 | 85.21 | 90.78 | 95.00 |
| Deep+LOMO [11] | 55.51 | 78.40 | 83.68 | 92.59 |
| Mirror [12] | 40.40 | 64.63 | 75.34 | 84.08 |

related to CAN-VGG16 method on VIPeR database is because the size of training set of VIPeR is so small.

### 4.3. Case 2: Background information within SCNCD

This experiment analyzes the accuracy performance of the background/foreground information within SCNCD (Sec. 3.3.2), before the employment of the post-ranking approach (Sec. 3.4). To this end, we set up the adopted framework to load only the following visual cues, $\mathbf{C}_5$ and $\mathbf{C}_6$ (detailed in Sec. 3.3.1, *i.e.*, without deep features), both without and with background/foreground information. Fig. 5 shows the CMC curves obtained for this experiment (for the first rank values). As it can be observed, the background/foreground information significantly improved the overall accuracy on the three evaluated datasets, being effective to remove the background noise. Yang et al. [6] obtained same conclusion when evaluating both representations (image-foreground and image-only) on VIPeR and PRID450s datasets. However, differently from their work, in which the evaluation was performed using only RGB information combined with the segmentation model proposed in [38] and the KISSME [9] metric learning, we adopted a more powerful segmentation strategy, as well as a different similarity function.

It can also be noticed from Fig. 5 that the proposed feature representation based on SCNCD outperformed obtained results (for the VIPeR dataset) reported in [5] (see Table 3).
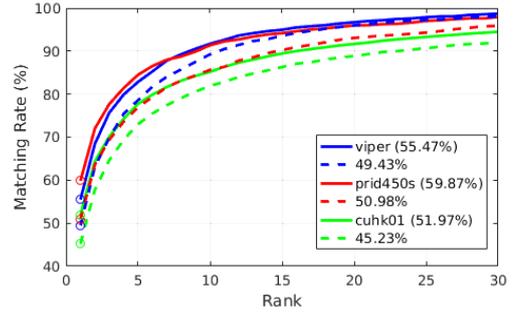


Figure 5: Accuracy performances based on SCNCD (*i.e.*, using only $\mathbf{C}_5$ and $\mathbf{C}_6$), with and without background/foreground information (solid and dashed lines, respectively). Top-1 rank values for each case are also provided.

### 4.4. Case 3: Complementary feature representations

This experiment evaluated the complementary features individually (before post-ranking). Each proposed feature representation was integrated as detailed in Sec. 3.3.5, and $s_i^{final}$ is adopted as the similarity function related to each representation $F_i$. Obtained results are shown in Fig. 6 in terms of top-1 rank recognition rate, from where we can make the following observations:

- All complementary features outperformed the baseline feature representation $F_0$.

- The benefits of including GOG feature into the baseline feature representation can be observed if we compare overall results obtained for the respective pair of feature representations $\langle F_0, F_1 \rangle$. Similarly, the benefits of including deep feature can be observed if we compare results for the pair $\langle F_0, F_2 \rangle$.

- The inclusion of both visual cues based on GOG and deep features, either on the baseline feature representation or on the proposed complementary features, can be highlighted if we compare overall obtained results for the respective pairs $\langle (F_1, F_2), F_3 \rangle$ and $\langle (F_4, F_5), F_6 \rangle$, as well as in relation to their respective simplified versions $\langle (F_7, F_8), F_9 \rangle$ and $\langle (F_{10}, F_{11}), F_{12} \rangle$.

- The simplified versions of the complementary features $F_9$ and $F_{12}$ obtained better accuracy than their respective complete representations ($F_3$ to $F_6$), indicating that the proposed simplification still has strong discriminative power for person re-identification applications, while requiring less computation resources.

- $F_6$ and $F_{12}$, which exploits SCNCD (with background/foreground information), GOG and deep features, obtained the best overall accuracy performance in the three adopted datasets.

The previously mentioned observations indicate that the proposed complementary feature representations have
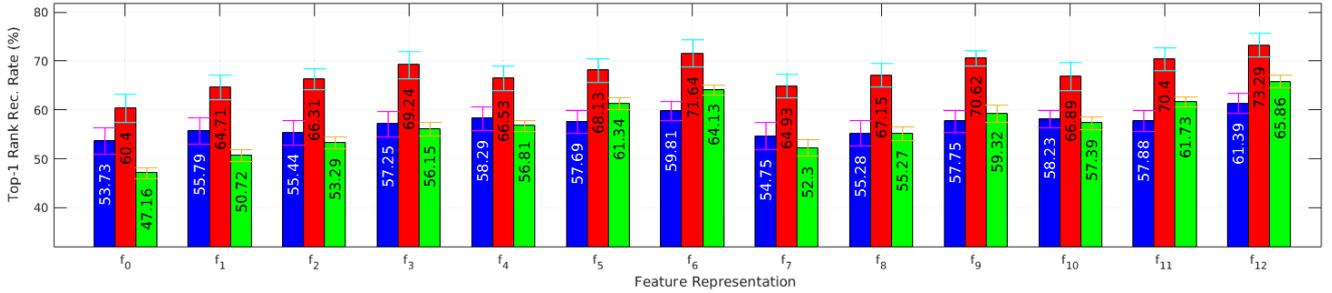
Figure 6: Accuracy performance obtained for each feature representation $F_i$ (before post-ranking), for the VIPeR (blue), PRID450s (red) and CUHK01 (green) datasets.

strong discriminative power in person re-identification applications, mainly when combined through ranking aggregation, as shown in Sec. 4.2. As it will be described in Sec. 4.5, accuracies obtained by such complementary features can also be improved by the post-ranking approach. In addition, different integration strategies (from those described in Sec. 3.3.5) were also evaluated in other experiments (*e.g.*, the integration of all features, $C_1$ to $C_9$, using the simplified and complete representations), however, no significant accuracy performance improvements were observed.

### 4.5. Case 4: Post-ranking analysis

In this section we analyze the results obtained by the DCIA method applied to the proposed framework. Our experiments were performed without visual expansion (employed in [10]), which synthesizes the probe into the gallery feature space aiming to reduce feature inconsistency. We avoided using this procedure because of its high computational requirements and, as commented by the authors, obtained accuracies were improved by less than 1%. Figure 7 shows results (averaged per database) for each complementary feature, before and after post-ranking. As we can observe, the post-raking approach improved overall results.
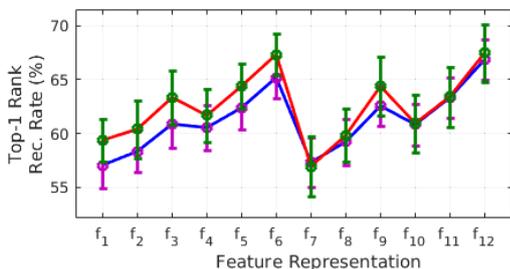


Figure 7: Average accuracy performance obtained for each feature representation, with (red) and without (blue) post-ranking.

As related in [10], one limitation of the DCIA method is that sometimes the true match is not included in the content set. In this case, it will not be re-ranked. In addition, some images might move to higher rank positions after post-ranking. However, in general, it improved more

than deteriorated final results. Fig. 8 shows the obtained improvement by the post-ranking approach (after rank-aggregation) on the three employed databases.
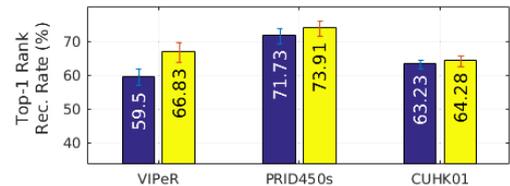


Figure 8: Average accuracy performance (after rank-aggregation) without/with post-ranking, respectively. Values obtained using the best-12 aggregation strategy described in Sec. 4.6.

Table 4 shows statistics related to the post-ranking approach (before ranking aggregation), considering 10 different runs and all complementary features $\{F_1, ..., F_{12}\}$. First, values were averaged in relation to each feature representation and number of runs. Then, the average value per database was computed.

Table 4: Post-ranking statistics: mean average and standard deviation (%) of: (i) probe images included in the content set, (ii) improved results, (iii) from the improved results, the ones that were moved to top-1 rank position, (iv) unchanged ranks and (v) images that were moved to higher rank positions (worsen).

|     | VIPeR | PRID450s | CUHK01 |
|-----|-------|----------|--------|
| i   | 76.3 ±2.8 | 81.2 ±1.7 | 71.3 ±1.4 |
| ii  | 11.4 ±2.4 | 7.0 ±1.8 | 6.4 ±1.1 |
| iii | 78.1 ±6.8 | 84.0 ±9.1 | 81.3 ±7.3 |
| iv  | 80.5 ±2.9 | 86.3 ±2.5 | 85.7 ±1.8 |
| v   | 8.1 ±2.1 | 6.6 ±1.8 | 7.9 ±1.5 |

From Table 4, we can observe that a high percentage of images which rank position was improved (ii), were moved to top-1 rank position (iii). On the other hand, few images were undesired moved to higher rank positions (v). However, post-ranking improved overall results.

### 4.6. Case 5: best-n tunning rank-aggregation strategy

In this experiment, the train set of each dataset was divided into 50% train and 50% validation. Fig. 9 shows obtained results on the validation set (averaged per database)

for each complementary feature (with post-ranking). Table 5 shows obtained results, on the test set, from the aggregation of the best-$n$ feature representations, for $n = 2$ to 12. As we can observe from Table 5, the inclusion of additional feature representations do not always increase accuracy performance, as well as there is not an overall feature representation that fits all databases.
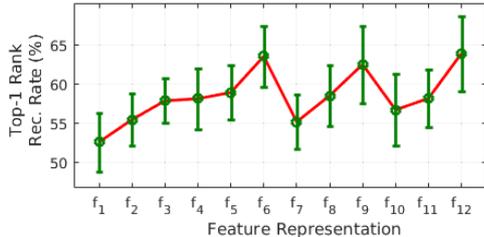


Figure 9: Average accuracy performance obtained for each feature representation on the validation set.

Table 5: Top-1 rank recognition rate (%) after ranking-aggregation, using the "best-$n$" feature representations.

| best-$n$ | VIPeR | PRID450s | CUHK01 |
|---|---|---|---|
| 2 | 64.75 | 73.87 | 64.90 |
| 3 | 66.99 | **75.64** | **66.91** |
| 4 | 66.87 | 75.24 | 66.75 |
| 5 | 66.83 | 74.56 | 66.30 |
| 6 | 66.77 | 74.31 | 66.69 |
| 7 | 66.33 | 74.18 | 66.34 |
| 8 | 66.83 | 74.98 | 65.05 |
| 9 | 67.15 | 75.07 | 66.13 |
| 10 | **67.21** | 74.40 | 65.49 |
| 11 | 66.77 | 73.95 | 65.10 |
| 12 | 66.83 | 73.91 | 64.28 |

### 4.7. Computational cost

We adapted the MATLAB implementation provided in [5] to consider the proposed complementary features. Computational costs shown in Table 6 were obtained using the VIPeR dataset. The complete representations $\{F_3, F_6\}$, which explore different visual cues, as well as their respective simplified versions $\{F_9, F_{12}\}$ were analyzed[8].

Average computational time to run the ranking-aggregation, using the best-7 strategy (described in Sec. 4.6) was 3.9m ±0.4.

### 5. Conclusion

In this work we exploited different feature representations, combined with a post-ranking and ranking aggregation strategies, to advance the state-of-the-art in person re-identification. Our model was built on a framework combining similarity learning metric with spatial constraints.

---

[8]Using a 2.30GHz Intel Core i7 CPU and 8Gb of memory, without considering I/O procedures and image resize operations.

Table 6: Average computational cost obtained from 10 runs, to process (train and test) the whole VIPeR database (316 train and 316 test images), taking into account different features representations.

| | Total | Post-rank | Test |
|---|---|---|---|
| $F_3$ | 16.9m ±4.2 | 6.1m ±0.4 | 4.17s ±1.5 |
| $F_6$ | 5.0m ±0.7 | 3.3m ±0.3 | 2.19s ±0.6 |
| $F_9$ | 7.3m ±1.9 | 4.1m ±0.3 | 2.93s ±1.1 |
| $F_{12}$ | 3.7m ±0.4 | 2.5m ±0.4 | 1.91s ±0.5 |

The proposed complementary features demonstrated to have strong discriminative power, as well as to complement each other even when the simplified versions are employed. Different feature representations were analyzed individually and incrementally. The post-ranking approach demonstrated to be a powerful tool in person re-identification tasks, being able to improve initial results which could be further enhanced by the ranking-aggregation strategy. We show that handcrafted and deep features fusion enhance re-identification performance especially in domains where there is a reduced amount of available data. As a result, we improved the top-1 rank recognition by 2.43% and 2.66% on VIPeR and PRID450s datasets, respectively, as well as obtained competitive results on the CUHK01 database.

### Acknowledgements

### References

[1] R. Vezzani, D. Baltieri, R. Cucchiara, People reidentification in surveillance and forensics: A survey, ACM Computing Surveys 46 (2) (2013) 29:1–29:37.

[2] A. Bedagkar-Gala, S. K. Shah, A survey of approaches and trends in person re-identification, Image and Vision Computing 32 (4) (2014) 270 – 286.

[3] S. Gong, M. Cristani, S. Yan, C. L. (Eds.), Person Re-Identification, Advances in Computer Vision and Pattern Recognition. Springer, 2014.

[4] D. Cheng, Y. Gong, S. Zhou, J. Wang, N. Zheng, Person re-identification by multi-channel parts-based cnn with improved triplet loss function, in: CVPR, 2016, pp. 1335–1344.

[5] D. Chen, Z. Yuan, B. Chen, N. Zheng, Similarity learning with spatial constraints for person re-identification, in: CVPR, 2016, pp. 1268–1277.

[6] Y. Yang, J. Yang, J. Yan, S. Liao, D. Yi, S. Z. Li, Salient color names for person re-identification, in: ECCV, 2014, pp. 536–551.

[7] D. Chen, Z. Yuan, G. Hua, N. Zheng, J. Wang, Similarity learning on an explicit polynomial kernel feature map for person re-identification, in: CVPR, 2015, pp. 1565–1573.

[8] L. Bazzani, M. Cristani, V. Murino, Symmetry-driven accumulation of local features for human characterization and re-identification, Computer Vision and Image Understanding 117 (2) (2013) 130–144.

[9] M. Kstinger, M. Hirzer, P. Wohlhart, P. M. Roth, H. Bischof, Large scale metric learning from equivalence constraints, in: CVPR, 2012, pp. 2288–2295.

[10] J. García, N. Martinel, A. Gardel, I. Bravo, G. L. Foresti, C. Micheloni, Discriminant context information analysis for post-ranking person re-identification, IEEE Transactions on Image Processing 26 (4) (2017) 1650–1665.

[11] S. Wu, Y. C. Chen, X. Li, A. C. Wu, J. J. You, W. S. Zheng, An enhanced deep feature representation for person re-identification, in: IEEE Winter Conf. on Applications of Computer Vision (WACV), 2016.

[12] Y.-C. Chen, W.-S. Zheng, J. Lai, Mirror representation for modeling view-specific transform in person re-identification, in: International Joint Conf. on Artificial Intelligence, 2015, pp. 3402–3408.

[13] T. Xiao, H. Li, W. Ouyang, X. Wang, Learning deep feature representations with domain guided dropout for person re-identification, in: CVPR, 2016, pp. 1249–1258.

[14] S. Paisitkriangkrai, C. Shen, A. van den Hengel, Learning to rank in person re-identification with metric ensembles, in: CVPR, 2015, pp. 1846–1855.

[15] W. Li, R. Zhao, T. Xiao, X. Wang, Deepreid: Deep filter pairing neural network for person re-identification, in: CVPR, 2014, pp. 152–159.

[16] D. Gray, S. Brennan, H. Tao, Evaluating appearance models for recognition, reacquisition, and tracking, in: IEEE Int. Workshop on Performance Evaluation for Tracking and Surveillance, 2007.

[17] R. F. de Carvalho Prates, W. R. Schwartz, CBRA: Color-based ranking aggregation for person re-identification, in: IEEE International Conference on Image Processing (ICIP), 2015, pp. 1975–1979.

[18] P. M. Roth, M. Hirzer, M. Koestinger, C. Beleznai, H. Bischof, Mahalanobis distance learning for person re-identification, in: S. Gong, M. Cristani, S. Yan, C. C. Loy (Eds.), Person Re-Identification, Springer, London, United Kingdom, 2014, pp. 247–267.

[19] W. Li, R. Zhao, X. Wang, Human reidentification with transferred metric learning, in: ACCV, 2012, pp. 31–44.

[20] N. Dalal, B. Triggs, Histograms of oriented gradients for human detection, in: CVPR, 2005, pp. 886–893.

[21] S. Liao, G. Zhao, V. Kellokumpu, M. Pietikinen, S. Z. Li, Modeling pixel process with scale invariant local patterns for background subtraction in complex scenes, in: CVPR, 2010, pp. 1301–1306.

[22] P. Luo, X. Wang, X. Tang, Pedestrian parsing via deep decompositional network, in: ICCV, 2013, pp. 2648–2655.

[23] T. Matsukawa, T. Okabe, E. Suzuki, Y. Sato, Hierarchical gaussian descriptor for person re-identification, in: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 1363–1372.

[24] S. Liao, Y. Hu, X. Zhu, S. Z. Li, Person re-identification by local maximal occurrence representation and metric learning, in: CVPR, 2015, pp. 2197–2206.

[25] B. Mirmahboub, M. L. Mekhalfi, V. Murino, Distance penalization and fusion for person re-identification, in: 2017 IEEE Winter Conference on Applications of Computer Vision (WACV), 2017, pp. 1306–1314.

[26] J. M. Stuart, E. Segal, D. Koller, S. K. Kim, A gene-coexpression network for global discovery of conserved genetic modules, Science 302 (5643) (2003) 249 – 255.

[27] J. García, N. Martinel, C. Micheloni, A. Gardel, Person re-identification ranking optimisation by discriminant context information analysis, in: 2015 IEEE International Conference on Computer Vision (ICCV), 2015, pp. 1305–1313.

[28] S. Bai, X. Bai, Q. Tian, Scalable person re-identification on supervised smoothed manifold, in: CVPR, 2017, pp. 2530–2539.

[29] Y. Zhang, B. Li, H. Lu, A. Irie, X. Ruan, Sample-specific svm learning for person re-identification, in: CVPR, 2016, pp. 1278–1287.

[30] E. Ahmed, M. Jones, T. K. Marks, An improved deep learning architecture for person re-identification, in: CVPR, 2015, pp. 3908–3916.

[31] H. Liu, J. Feng, M. Qi, J. Jiang, S. Yan, End-to-end comparative attention networks for person re-identification, IEEE Transactions on Image Processing 26 (7) (2017) 3492–3506.

[32] Y. Du, H. Ai, S. Lao, Evaluation of color spaces for person re-identification, in: ICPR, 2012, pp. 1371–1374.

[33] D. Gray, H. Tao, Viewpoint invariant pedestrian recognition with an ensemble of localized features, in: ECCV, 2008, pp. 262–275.

[34] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, Q. Tian, Scalable person re-identification: A benchmark, in: ICCV, 2015.

[35] R. Kolde, S. Laur, P. Adler, J. Vilo, Robust rank aggregation for gene list integration and meta-analysis, Bioinformatics 28 (4) (2012) 573.

[36] Y. C. Chen, X. Zhu, W. S. Zheng, J. H. Lai, Person re-identification by camera correlation aware feature augmentation, IEEE Transactions on Pattern Analysis and Machine Intelligence PP (99) (2017) 1–14.

[37] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, X. Tang, Spindle net: Person re-identification with human body region guided feature decomposition and fusion, in: CVPR, 2017, pp. 1077–1085.

[38] N. Jojic, A. Perina, M. Cristani, V. Murino, B. Frey, Stel component analysis: Modeling spatial correlations in image class structure, in: CVPR, 2009, pp. 2044–2051.