

Multi-shot SURF-based Person Re-identification via sparse representation

Mohamed Ibn Khedher¹, Mounîm A. El Yacoubi², Bernadette Dorizzi¹

¹ Institut Mines-Telecom / Telecom SudParis : UMR5157, SAMOVAR, Evry, France

² Institut Mines-Telecom / Telecom SudParis : CEA Saclay Nano-Innov, 91191 Gif sur Yvette Cedex, France

{Mohamed.ibn.khedher, Mounim.El.Yacoubi, bernadette.dorizzi}@telecom-sudparis.eu

Abstract

We present in this paper a multi-shot human re-identification system from video sequences based on SURF matching. Our contribution is about the matching step which is crucial. In this context, we propose a new method of SURF matching via sparse representation. Each SURF Interest Point in the test sequence is represented by a sparse representation of SURFs points in the reference dataset. For efficiency purposes, a dynamic dictionary is selected for each SURF from this dataset through KD-Tree Neighborhood search. Then a majority vote rule is applied to classify the test sequence. This approach is evaluated on two public datasets : PRID-2011 and CAVIAR4REID. The experimental results show that our approach compares favorably with and outperforms current state-of-the-art on the two datasets by 1% to 7%.

1. Introduction

Person re-identification is an important task in visual surveillance due to its diverse applications (tracking criminals over multiple cameras, finding lost people, etc.) in different places (shopping centers, hospitals, streets, etc.). In a camera network, and given two cameras possibly having different scene views, if a person leaves the view of one camera and reappears in the other, the re-identification system must be able to re-identify him/her and continue monitoring [23].

The re-identification task is challenging since it can be sensitive to several factors such as varying appearances of a person across the camera network. In fact, a person may look different due to changes in camera characteristics, camera viewpoints, light conditions, poses, partial occlusions, etc.

Methods of re-identification can be single-shot-based [17]) or multiple-shot-based [16, 18]). In order to compare re-identification system performances, several datasets are available. For the single-shot method, we find VIPeR [17], PRID-2011 single-shot version [20]. For the multi-shot method, we find CAVIAR4REID [12], ETHZ [28] and PRID-

2011 multi-shot version [20].

Overall, methods of re-identification can be classified into two main approaches : global approaches and local approaches.

The former often exploit appearance features that include color and texture and are commonly used in the state-of-the-art. For instance, Bak et al. [4] build an appearance model based on Haar-like features and Dominant Color Descriptors (DCD). Then, in order to obtain the most invariant and discriminative signature, an AdaBoost scheme is applied to descriptors. In [14], several features are combined to model three complementary aspects of the human appearance : the overall chromatic content, the spatial arrangement of colors into stable regions, and the presence of recurrent local motifs with high entropy. Hirzer et al. [19] compute a descriptive appearance representation encoding the vertical color structure of a person and estimate the transition between two cameras using a pair-wise estimated metric to improve the classification results. In [27], an appearance representation based on the Major Color Spectrum Histogram Representation (MCSHR) is considered to cope with the typical illumination changes. In [8], Chromatic and Epitomic analyses are proposed to model the human appearance. They incorporate complementary global and local statistical descriptions of the human appearance. Bedagkar-Gala et al. [9] propose an adaptive part-based spatio-temporal model that characterizes person's appearance using color and facial features.

For local approaches, on the other hand, a person silhouette is represented by several locally detected points or patches. For instance, Gheissari et al. [16] extract spatio-temporal interest points described by color and structural information, then try to fit a triangulated graph to each person to cope with pose variation. In [22], Implicit Shape Models (ISM) based on SIFT descriptors are used to capture the shape properties of a person. Interest point-based local descriptors like PCA-SIFT and SURF are used in [2, 18] to capture the local appearance variations. In [13], SURF and color descriptors are combined to improve matching. In a comparative study of state-of-the-art [6], interest point

detectors and descriptors are evaluated for the task of re-identification. Finally, Bingpeng et al. [26] propose a new local descriptor for person re-identification encoded by Fisher Vectors.

One of the main observations regarding state-of-the-art re-identification methods is that there is no approach systematically outperforming the others : each approach has its own strengths and limitations. In this work, we consider a re-identification local approach based on Interest Points because of their relative robustness towards camera view angle change. State-of-the-art shows that many interest points' detectors and descriptors are used. Interest points are employed in different fields such as object recognition [25], face recognition [10] and action recognition [3]. Each method is characterized by its own description, invariance criteria and running time. In the state-of-the-art, no method systematically outperforms others since evaluations of different Interest Points show conflicting results. We have carried out some initial investigation regarding SIFT and SURF robustness on CASIA-B database [21]; they showed that SURF outperforms SIFT. Hence, in this work, SURF [7] is used. To overcome the instability of these points, we follow the multi-shot re-identification approach using all images in order to increase the reproducibility of interest points between two similar video scenes.

Several works based on the multi-shot approach and Interest Points Pair Matching have been considered [18, 13, 23]. These previous works share the same matching step where each test Interest Point is matched to the closest reference Interest Point. However, they differ in the matching pair selection (Filtering step). In [18], an empirically preset number of best matched points between query and reference is chosen and a majority vote scheme is used to validate a re-identification decision. In [13], a reference point p_0 is matched to a test point p_1 , if $d(p_1, p_0) < s d(p_1, p_i)$ ($\forall p_i \in$ Reference, where s is a preset threshold, $s < 1$ and $d(., .)$ is the Euclidian Distance). In [23], the empirical estimation of a threshold is avoided and an automatic statistical method of acceptance and rejection of SURFs correspondence based on the likelihood ratio of two GMM is proposed. Our work is different from [18, 13, 23] in the nature of Interest Points matching itself. We can notice that in all the previous methods, Interest Point matching consists of determining for each test Interest Point the closest reference one independently of the nature of video sequences. In fact, in real conditions where the environment is uncontrolled and many people can walk within the camera view, SURFs are much noisier and more ambiguous. Thus, the information of the closest reference SURF is insufficient to identify the test SURF.

Our idea consists of representing each test SURF as the sparsest linear combination of reference SURF points, and then exploiting the latter to identify the test SURF. Sparse

representation of signals has been studied since two decades but only recently it became popular after its application for face recognition [31]. Since then, it also has been used in other computer vision fields such as gait recognition, speech recognition and person re-identification. Our approach is different from [31] in 2 points. First, Sparse representation in [31] is global (Face representation as a whole), while ours is local (Sparse representation for each SURF in the silhouette). Second, in [31], for each sparse representation, the dictionary consists of all reference samples, while in our approach we use a reduced and dynamic dictionary of few selected reference SURFs as will be explained in section 3. Our approach is scalable to large datasets. It is different from [30] where one global sparse representation for each silhouette is found. Our contribution consists of computing for each test SURF a local sparse representation, independently of the other SURFs. For each representation, a dynamic dictionary of small size is generated. This dictionary is dynamic in the sense that it changes for each SURF and is selected based on a preset number of closest reference SURFs generated by a KD-Tree Neighborhood search.

This paper is structured as follows. In section 2, we describe the principle of sparse representation in the context of SURF matching. In section 3, we present the major stages of our re-identification system. The experimental results of our approach are given in section 4. Conclusions and perspectives are finally presented.

2. SURF-based Sparse Representation

The main idea of sparse representation is to find a representation of a signal involving the smallest number of elements of a preselected dictionary. Given a signal $y \in \mathbb{R}^D$ and a dictionary $\Phi \in \mathbb{R}^{D \times K}$ ($D \ll K$), there is an infinite number of solutions α verifying :

$$y = \Phi \alpha \quad (1)$$

The objective of sparse representation is to find the sparsest solution α_s of Equation 1. In the context of SURFs matching, let us consider a query SURF y and a set of reference SURFs associated with M identities (persons). First, the reference dataset is arranged into a matrix (called dictionary), which is built using reference SURFs : $\{S_{i,j}\} \in \mathbb{R}^D$, $i=1 \dots M$, $j=1 \dots k_i$, where k_i denotes the number of reference SURFs for the i -th identity, and $K=k_1+k_2+\dots+k_M$ denotes the number of SURFs in the reference dataset. The k_i reference SURFs of the i -th identity constitute the columns of the matrix Φ_i :

$$\Phi_i = [S_{i,1}; S_{i,2}; \dots; S_{i,k_i}] \quad (2)$$

All K SURFs from the reference dataset are combined to form the matrix Φ :

$$\Phi = [\Phi_1; \Phi_2; \dots; \Phi_M] = [S_{1,1}; S_{1,2}; \dots; S_{M,k_M}] \quad (3)$$

Hence, we represent the query SURF y as a linear combination of all the reference SURFs :

$$y = \Phi \alpha_s = [\Phi_1; \Phi_2; \dots; \Phi_M] \alpha_s \quad (4)$$

where α_s is a sparse coefficient vector and ideally can be written as the following :

$$\alpha_s = [0, \dots, 0, \alpha_{i,1}, \alpha_{i,2}, \dots, \alpha_{i,k_i}, 0, \dots, 0]^T \quad (5)$$

In this ideal case, α_s has non-zero entries associated only with the i -th identity corresponding to the real identity of y . In noisy conditions, however, coefficients associated with other identities may be non null.

The sparse representation problem consists of finding the sparsest solution α_s of Equation 4. Mathematically, the problem can be written as l_0 -norm minimization :

$$\alpha_s = \min_{\alpha} \|\alpha\|_0 \quad \text{subject to} \quad y = \Phi \alpha \quad (6)$$

To find the sparsest solution α_s without NP-hard complexity, it is sufficient to solve an l_1 -norm minimization problem :

$$\alpha_s = \min_{\alpha} \|\alpha\|_1 \quad \text{subject to} \quad y = \Phi \alpha \quad (7)$$

To take into account noise, the problem of Equation 7 can be generalized to the LASSO (Least Absolute Shrinkage and Selection Operator) [29] formulation as following :

$$\alpha_s = \min_{\alpha} \|\Phi \alpha - y\|_2^2 \quad \text{subject to} \quad \|\alpha\|_1 \leq \epsilon \quad (8)$$

where ϵ is a preset threshold. Another equivalent formulation to LASSO uses a scalar regularization parameter to adjust the tradeoff between sparsity and error reconstruction :

$$\alpha_s = \min_{\alpha} (\|\Phi \alpha - y\|_2^2 + \lambda \|\alpha\|_1) \quad (9)$$

After calculating the sparsest solution, the non-zero coefficients of α_s can be used to determine the identity of the query SURF y .

3. The person re-identification system

Our approach basically consists of three stages : 1) Feature extraction using SURF, 2) SURF identification via sparse representation and 3) Person re-identification based on majority vote rule. Figure 1 shows the flowchart of our approach.

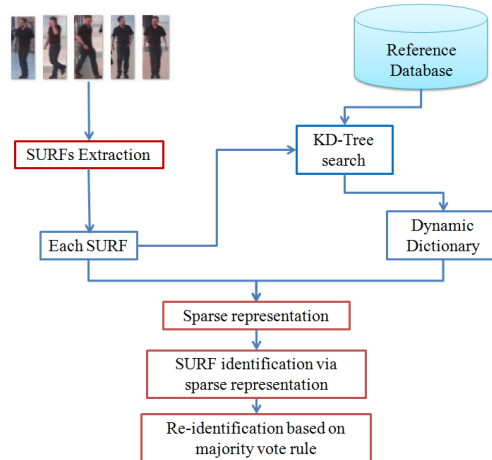


FIGURE 1. Re-identification stages

3.1. Feature extraction

Several interest point detectors and descriptor are found in the literature like SIFT and SURF. SURF [7] is used in this work due to its fast computation and relatively scale and rotation invariance. SURF operates in two main stages, namely the detector and the descriptor stages. SURF uses the determinant of the Hessian matrix to detect the interest points. The detector analyzes an image and returns a list of centers of interest points. The SURF descriptor captures information from the region around the interest point. This description is built from local intensity differences. In fact, a SURF descriptor of dimension 64 is computed as a sum of local intensity differences within a 4x4 grid around the interest point. These intensity differences are calculated as responses to first-order Haar-Wavelets. For illumination invariance, the descriptor is normalized to unit length [23]. Figure 2 shows the detected SURF points within a sample from the used dataset.

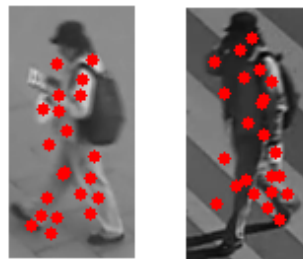


FIGURE 2. Feature extraction

3.2. SURF identification via sparse representation

The re-identification system consists of classifying each test SURF independently. To identify one SURF, three steps are applied : 1) Dictionary construction, 2) Sparse representation and 3) Identity assignment.

3.2.1 Dictionary construction

As explained above, generally, the dictionary consists of all elements of the reference dataset. Thus for typical re-identification datasets, we would have to consider dictionaries of millions of reference SURFs. The computation time for finding a sparse representation for one SURF will then be huge. In this work, for each test SURF a dynamic and reduced dictionary A is chosen consisting of the N closest SURFs from the reference dataset. The dimension of A is $D \times N$ where each column represents a SURF description of dimension $D=64$ and N is empirically set to 200. To accelerate the search of Nearest Neighbors, a KD-tree is used [15].

3.2.2 Sparse representation

In our experiments, the popular sparse representation algorithm LASSO is chosen because it takes into account the fact that SURFs are noisy. This algorithm takes as input one test SURF and the corresponding dictionary dynamically built and outputs a vector with most coefficients equal to zero. Figure 3 shows the inputs and outputs of the LASSO algorithm.



FIGURE 3. Inputs and outputs of LASSO

3.2.3 Identity assignment

After calculating the sparse representation of SURF, the question becomes : how can we use non-zero coefficients of the previous representation to assign an identity to a test SURF ? To perform such an assignment, a residual is calculated for each identity i having at least one non-zero coefficient in the following manner :

Denote L the number of identities having at least one non-zero coefficient in α_s . Let y be a SURF description and A the corresponding dictionary.

We compute first $x_i : i=1 \dots L$, as following :

$$x_i = [0, \dots, 0, \alpha_{i,1}, \alpha_{i,2}, \dots, \alpha_{i,k_i}, 0, \dots, 0] \quad (10)$$

x_i is a coefficient vector obtained from α_s with all elements zero except those associated with the identity i . The closest identity for the test (query) SURF satisfies the following equation :

$$j = \arg \min_i (\|y - Ax_i\|_2^2) \quad (11)$$

j corresponds to the identity minimizing the reconstruction residual of y . Thus, the identity of the test SURF y is assigned to j .

3.3. Human re-identification

Human re-identification basically consists of two stages. 1) SURF identification via sparse representation and 2) majority vote decision rule. Given a test sequence and a reference dataset, the objective is to assess whether a sequence from the same person as the test sequence is within the reference dataset. In the first step, each SURF from the test sequence is classified into one identity from the reference dataset via sparse representation as explained in section 3.2. In the second stage, the found reference identities are submitted to the majority vote decision rule. For each test SURF, a vote is added to the person associated with the reference selected identity. The person obtaining the majority of votes is claimed as the re-identified person.

4. Experimental results

We evaluated our approach on two public datasets : CAVIAR4REID [12] and PRID-2011 multi-shot version [20]. Results are shown in terms of the Cumulative Matching Characteristic (CMC) curve associated with the identification rate as commonly used in the literature.

4.1. Results on CAVIAR4REID

CAVIAR4REID [12] has been extracted from the CAVIAR database [11]. The recorded videos were captured from two different cameras in an indoor shopping center in Lisbon. The pedestrians' images have been cropped using the provided ground truth. From the 72 different individuals identified (with image sizes varying from 17x39 pixels to 72x144 pixels), 50 people are captured by both views and 22 from only one camera. For each pedestrian, 10 images from each camera view are selected, maximizing the variance with respect to resolution changes, light conditions, occlusions, and pose changes (see samples in Figure 4).



FIGURE 4. Some images of CAVIAR4REID, each pair portraying the same individual

The CMC curve obtained by our approach is shown in Figure 6. Table 1 shows different methods performances (identification rate at rank 1) found in the state-of-the-art. The approaches based on appearance features (SDALF) [24] and MRGC [5] using essentially color descriptors achieve 10% of correct identification. The approach presented in [23] based on SURF matching and probabilistic filtering achieves 16% of correct identification. Our approach

achieves an identification rate of 18%, which is slightly higher than the one obtained in [12]¹ and based on spatial-temporal color features. CMC curves sources of [24, 12] on CAVIAR4REID are not available to reproduce along with ours in the same figure.

TABLE 1. Results comparison on CAVIAR4REID

Approach	Re-identification rate (%)
Authors of [5] appear in [24]	10
[24]	10
[23]	16
[12]	17
Our approach	18

4.2. Results on PRID-2011

The dataset PRID-2011² (multi-shot version) [20] was created in 2011 by the Austrian Institute of Technology. The movies were obtained from two cameras (A and B) located on the street (Figure 5). 385 people were filmed by camera-A and 749 people were filmed by Camera-B (200 are common to the two cameras). The evaluation consists of searching the common 200 people filmed by Camera-A in the gallery set (Camera-B) of 749 people.



FIGURE 5. Samples images from PRID-2011. Upper and lower rows correspond to different camera views

Figure 6 shows the obtained CMC (From rank 1 to rank 10) on the PRID-2011 dataset compared to the state-of-the-art. Table 2 shows different methods performances (identification rate at rank 1). The approach of [23] based on SURF matching and probabilistic filtering gives 22% of correct identification. [20] achieves 19.18%³ of correct identification. Authors of [20] combine an appearance descriptor based on a set of region covariance and a discriminative model based on boosting feature selection. Our approach achieves 29% of correct identification.

1. The work in [12] was carried out by the team that produced CAVIAR4REID dataset.

2. This is actually a cleaned version [1].

3. Note that although the authors in [20] refer to the whole Prid-2011 dataset, their results are actually obtained on the cleaned version, as shown in the website [1].

TABLE 2. Results comparison on PRID-2011.

Approach	Re-identification rate (%)
[20]	19.18
[23]	22
Our approach	29

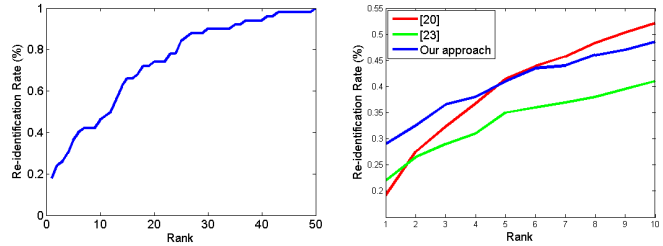


FIGURE 6. Left :CMC performance on CAVIAR4REID of our approach, Right : Comparison of CMC performances on PRID-2011

4.3. Comparison of results

Note that our approach and [23], both are based on SURFs matching. On PRID-2011, our approach outperforms [23] at all ranks and achieves an improvement of 7% in the rate of re-identification at rank 1. This improvement is significant given the large size of the dataset and proves that the sparse representation can provide richer information for decision making than [23] and other interest point matching methods like [18, 13]. Compared to [20], our approach becomes less efficient starting from rank 6. However, our approach outperforms [20] from rank 1 to rank 5 and achieves an improvement of 9.82 % in the rate of re-identification at rank 1. This improvement is also significant and shows the power of sparse representation compared to [20] that is the combination of two methods.

Sparse representation is better again on CAVIAR4REID. The improvement is small compared to the one obtained on PRID-2011, maybe because the database is small and few images only are available per person on the contrary of PRID-2011.

5. Conclusion

This paper has studied the performance of a multi-shot human re-identification system based on SURF matching. It proposed a new method of SURF matching via sparse representation that consists of representing each test SURF as the sparsest linear combination of reference SURFs. For efficiency, a dynamic dictionary is selected based on a preset number of closest reference SURFs obtained by KD-Tree Neighborhood search. The results obtained in our experiments show the relative power of sparse representation to match each Interest Point in noisy and ambiguous conditions that are inherent to real video sequences. Our approach compares favorably with the state of the art mainly on the

large dataset Prid-2011.

In the future, we will focus on other exploitations of the non-zero coefficients of the sparse representation in order to assign an identity to a SURF test. On the other hand, we will investigate the cooperation of re-identification methods based on color, Interest Points and/or other methods depending on the complexity of scene.

Références

- [1] <http://lrs.icg.tugraz.at/datasets/prid/index.php>, 2011.
- [2] C. Arth, C. Leistner, and H. Bischof. Object reacquisition and tracking in large-scale smart camera networks, 2007.
- [3] H. K. S. I. Atiqur Rahman Ahad, J. Tan. Surf- and optical flow-based action recognition with outlier management. In *IEEE Computer Vision and Pattern Recognition (CVPR) workshop on Gesture Recognition*, 2011.
- [4] S. Bak, E. Corvee, F. Bremond, and M. Thonnat. Person re-identification using haar-based and dcd-based signature. *Advanced Video and Signal Based Surveillance, IEEE Conference on*, 0 :1–8, 2010.
- [5] S. Bak, E. Corvee, F. Bremond, and M. Thonnat. Person re-identification using spatial covariance regions of human body parts. *Advanced Video and Signal Based Surveillance, IEEE Conference on*, 0 :435–440, 2010.
- [6] M. Bauml and R. Stiefelhagen. Evaluation of local features for person re-identification in image sequences. *Advanced Video and Signal Based Surveillance, IEEE Conference on*, 0 :291–296, 2011.
- [7] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool. Speeded-up robust features (surf). *Comput. Vis. Image Underst.*, 110(3) :346–359, 2008.
- [8] L. Bazzani, M. Cristani, A. Perina, and V. Murino. Multiple-shot person re-identification by chromatic and epitomic analyses. *Pattern Recogn. Lett.*, 33(7) :898–903, 2012.
- [9] A. Bedagkar-Gala and S. K. Shah. Part-based spatio-temporal model for multi-person re-identification. *Pattern Recogn. Lett.*, 33(14) :1908–1915, 2012.
- [10] M. Bicego, A. Lagorio, E. Grosso, and M. Tistarelli. On the use of sift features for face authentication. In *In : Conf. on Computer Vision and Pattern Recognition Workshop (CV-PRW). (2006)*, page 35, 2006.
- [11] CAVIAR. <http://homepages.inf.ed.ac.uk/rbf/caviar/>, 2003.
- [12] D. S. Cheng, M. Cristani, M. Stoppa, L. Bazzani, and V. Murino. Custom pictorial structures for re-identification. In *British Machine Vision Conference (BMVC)*, 2011.
- [13] I. O. de Oliveira and J. L. de Souza Pio. People reidentification in a camera network. *Dependable, Autonomic and Secure Computing, IEEE International Symposium on*, 0 :461–466, 2009.
- [14] M. Farenzena, L. Bazzani, A. Perina, V. Murino, and M. Cristani. Person re-identification by symmetry-driven accumulation of local features. In *CVPR*, pages 2360–2367, 2010.
- [15] J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. *ACM Trans. Math. Softw.*, 3(3) :209–226, 1977.
- [16] N. Gheissari, T. B. Sebastian, and R. Hartley. Person reidentification using spatiotemporal appearance. In *Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition - Volume 2*, pages 1528–1535, 2006.
- [17] D. Gray and H. Tao. Viewpoint invariant pedestrian recognition with an ensemble of localized features. In *Proceedings of the 10th European Conference on Computer Vision : Part I*, pages 262–275, 2008.
- [18] O. Hamdoun. *Détection et ré-identification de piétons par points d'intérêt entre caméras disjointes*. PhD thesis, École Nationale Supérieure des Mines de Paris, 2010.
- [19] M. Hirzer, C. Beleznaï, M. Koestinger, P. M. Roth, and H. Bischof. Dense appearance modeling and efficient learning of camera transitions for person re-identification. In *Proc. IEEE Int'l Conf. on Image Processing*, 2012.
- [20] M. Hirzer, C. Beleznaï, P. M. Roth, and H. Bischof. Person re-identification by descriptive and discriminative classification. In *Proceedings of the 17th Scandinavian conference on Image analysis*, pages 91–102, 2011.
- [21] <http://www.cbsr.ia.ac.cn/english/Gait>
- [22] K. Jungling and M. Arens. View-invariant Person Re-identification with an Implicit Shape Model. In *2011 8th IEEE International Conference on Advanced Video and Signal-Based Surveillance (AVSS)*, page 6, 2011.
- [23] M. I. Khedher, M. A. El-Yacoubi, and B. Dorizzi. Probabilistic matching pair selection for surf-based person re-identification. In *BIOSIG*, pages 1–6, 2012.
- [24] V. M. L. Bazzani, M. Cristani. Symmetry-driven accumulation of local features for human characterization and re-identification. *Computer Vision and Image Understanding*, 2013.
- [25] D. G. Lowe. Local feature view clustering for 3d object recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 682–688, 2001.
- [26] B. Ma, Y. Su, and F. Jurie. Local descriptors encoded by fisher vectors for person re-identification. In *Proceedings of the 12th international conference on Computer Vision - Volume Part I*, pages 413–422, 2012.
- [27] C. Madden, E. D. Cheng, and M. Piccardi. Tracking people across disjoint camera views by an illumination-tolerant appearance representation. *Mach. Vision Appl.*, 18 :233–247, 2007.
- [28] W. R. Schwartz and L. S. Davis. Learning discriminative appearance-based models using partial least squares. In *SIB-GRAPI*, pages 322–329, 2009.
- [29] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society (Series B)*, 58 :267–288, 1996.
- [30] N. Truong Cong, C. Achard, and L. Khoudour. People re-identification by classification of silhouettes based on sparse representation. In *International Conference on Image Processing Theory, Tools and Applications*, pages 60–65, 2010.
- [31] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 31(2) :210–227, 2009.