

# FEATURE EXTRACTION TECHNIQUES FOR ABANDONED OBJECT CLASSIFICATION IN VIDEO SURVEILLANCE

*Ahmed Fawzi Otoom, Hatice Gunes, Massimo Piccardi*

Faculty of Information Technology, University of Technology, Sydney (UTS)  
Sydney, Australia

E-mail: {afaotoom, haticeg, massimo}@it.uts.edu.au

## ABSTRACT

We address the problem of abandoned object classification in video surveillance. Our aim is to determine (i) which feature extraction technique proves more useful for accurate object classification in a video surveillance context (scale invariant image transform (SIFT) keypoints vs. geometric primitive features), and (ii) how the resulting features affect classification accuracy and false positive rates for different classification schemes used. Objects are classified into four different categories: bag (s), person (s), trolley (s), and group (s) of people. Our experimental results show that the highest recognition accuracy and the lowest false alarm rate are achieved by building a classifier based on our proposed set of statistics of geometric primitives' features. Moreover, classification performance based on this set of features proves to be more invariant across different learning algorithms.

**Index Terms** — Abandoned object classification, video surveillance, statistics of geometric primitives, SIFT keypoints.

## 1. INTRODUCTION

Automatic recognition, description, classification and grouping of patterns have been identified as significant problems within the computer vision research community and have been tackled for decades. In recent years, there has been growing interest and effort in developing research approaches for recognizing objects in still images. The majority of these approaches focus on extracting local regions such as Difference of Gaussian (DoG) regions [6], saliency regions [5], or other types of local patches. A discriminative model for recognition is then built based on these features such as: constellation models [4], "bag of words" models [11], and others. Results of these approaches are promising for objects categorization. However, the extracted features depend largely on local regions, such as corners and textured patches, therefore recognize objects only from one viewpoint and might not be accurate for recognizing objects when the viewpoint changes (e.g. [4]).

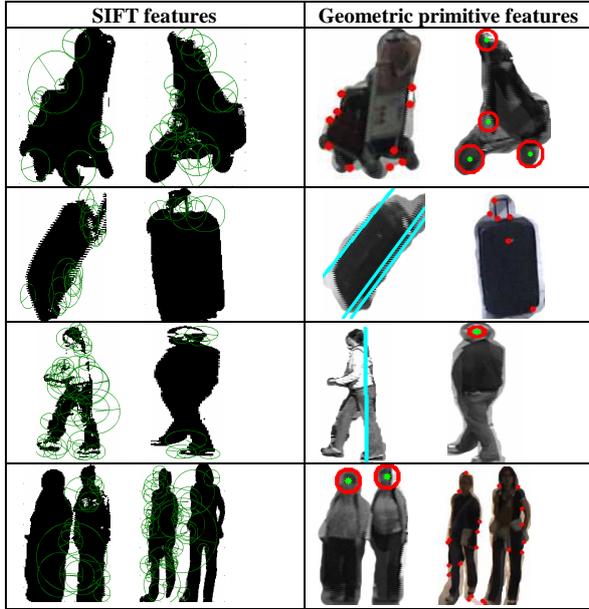
Object classification in video surveillance has also gained more attention recently. It aims to classify objects of interests into a number of predefined categories. Object categories are defined in advance depending on the environment where these objects are likely to be detected in the scene. Images of objects of interest are first analyzed in order to choose features that are simple yet

efficient to discriminate between the predetermined classes. Extracted features should be robust to various challenging conditions such as occlusion and change in viewpoint and illumination. In general, moving object recognition has gained more attention than abandoned object recognition [3, 10]. However, abandoned objects need to be detected and classified in an accurate way due to the fact that such objects may represent a high security threat. Efficient and accurate classification is needed in order to assess the potential danger they might cause prior to taking appropriate actions. Existing approaches for abandoned object recognition mainly depend on extracting a limited number of shape or appearance features [2, 9], resulting in a classifier that may not be capable of addressing the various challenges faced in a surveillance environment (e.g. [9]).

Within the rich body of literature on object and/or object class recognition, it is often stated that great attention should be paid to the definition of a discriminative feature set. There exist previous works for evaluating the performance of feature extraction techniques based on different local region descriptors and across a number of classifiers (e.g. [7]). However, there has been no attempt to compare local region features with statistics of geometric primitives' features in a visual surveillance context. Accordingly, in this paper, we aim to determine (i) which feature extraction technique proves more useful for accurate object classification in a video surveillance context (scale invariant image transform (SIFT) keypoints vs. geometric primitive features), and (ii) how the resulting features affect classification accuracy and false positive rates for different classification schemes used.

The work presented in this paper aims to become an integral part of a video surveillance system framework that is able to track multiple people and automatically detect abandoned objects for security of crowded areas such as a railway station or an airport terminal. Our work is based on the assumption that the abandoned object is already detected by a detector of "new stationary objects" in the scene; its location and size are also made available. A commercial off-the-shelf technology product (e.g., [14]) can be used for this task. We also assume that the area of interest is located within an airport or train station, and the objects of interest consist of trolley(s), bag(s), single person and group(s) of people. The problem at stake should not be confused with generic object classification, for which several methods exist suited to variable number and type of object classes ([4-7] and others): instead, given the high cost associated with misclassification errors in a

surveillance context, we aim to devise the most accurate feature extraction procedure possible given the categories of interest. The remainder of this paper is organized as follows: in Section 2 we introduce the feature extraction techniques. Classification learning methods and performance evaluation are described in Section 3. Experimental results and analysis are presented in Section 4. Finally, we draw our conclusions in Section 5.



**Figure 1.** Examples of features detected in a number of images: trolley (1<sup>st</sup> row), bag (2<sup>nd</sup>), person (3<sup>rd</sup>) and group of people (4<sup>th</sup>).

## 2. FEATURE EXTRACTION

The first step in any classification problem is feature extraction where features are extracted from images based on different image information. We apply three different approaches for extracting features. These approaches are based on SIFT keypoints and statistics of geometric primitives.

### 2.1 SIFT keypoints

SIFT (Scale-Invariant Feature Transform) keypoints are known to be invariant to rotation, scale, and translation, and are used to detect distinctive edges and textures in an image. Moreover, SIFT has empirically outperformed many other descriptors [7]. Because of the aforementioned reasons we choose to apply SIFT for the detection and description of local features (keypoints). Each keypoint is described with a 132-dimension vector: 128 spatial orientations, plus coordinates, scale, and rotation. After extracting SIFT keypoints from all images, we first apply dimensionality reduction and then we apply two different approaches for the final description of the features as illustrated in the following subsections. Figure 1 (left column) shows examples of SIFT keypoints detected in a number of images.

#### 2.1.1 Dimensionality reduction

After extracting SIFT keypoints, it is necessary to reduce the dimensionality in order to extract significant information and be

capable of training classifiers. We apply two popular dimensionality reduction techniques: principle component analysis (PCA) and linear discriminant analysis (LDA). From the initial analysis of the results, both techniques seem similar in their performance for the final classification results, with PCA slightly outperforming LDA. Therefore, we present PCA-based results. PCA is an orthogonal transformation of the coordinate system that describes the data. Given a set of  $M$  centered observations  $x_i \in R^N$ ,  $i = 1, \dots, m$ ,  $\sum_{i=1}^m x_i = 0$ , PCA finds the principle axes by diagonalizing the covariance matrix

$$C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T \quad (1)$$

To provide the diagonalization, the Eigenvalue equation  $\lambda v = Cv$  has to be solved where  $v$  is the Eigenvector matrix. The first few Eigenvectors are used as the basis vectors for the lower dimensional space. PCA aligns the data along the directions of the greatest variance. We keep only the eigenvectors corresponding to the highest eigenvalues, capturing 90% the variance within the data set. We thus reduce the dimensionality of the keypoint vectors down from 132 to 3. After applying PCA, we apply two approaches for the final description of the SIFT keypoints: majority rule approach and keypoint histograms approach.

#### 2.1.2 Approach 1: SIFT keypoints and majority rule

In this method, each keypoint in an image is classified independently and the final decision for the image class is the same class assigned to the majority of its keypoints. Let  $X$  be the class assigned to keypoint  $i$  in an image  $M$  and  $d(x | f_i)$  be the binary decision (0|1) for a keypoint  $i$  given feature vector  $f_i$ . Since  $X$  is one of four classes (person, group, bag, trolley), then  $d(x | f_i) = 1$  for only one class and 0 for all the others. For each image  $M$ , using the number of keypoints denoted as  $T$ , the multiple decisions are added up, for each class separately, as:

$$D(x | f_1, \dots, f_T) = \sum_{i=1}^T d(x | f_i). \quad \text{The final class assigned to}$$

$M$  will then be

$$x^* = \arg \max_x (D(x | f_1, \dots, f_T)) \quad (2)$$

#### 2.1.3 Approach 2: SIFT keypoint histograms

As our main goal is that of comparing feature extraction techniques, this approach was inspired by [1], except that we apply PCA instead of LDA for the feature reduction. We create a keypoint histogram for each image allowing the relationships between numbers and types of keypoints to be extrapolated and the information on the actual location discarded. Following this rationale, we first apply PCA to each keypoint, as explained in Section 2.1.1. Secondly, we choose a number of bins for each feature to be approximately proportional to the data variance within that feature. Eventually we use a histogram with 6, 4 and 2 bins for 1-3 features obtained from PCA. The resulting histograms are then fed into the classifiers for object classification.

### 2.2 Approach 3: Statistics of geometric primitives

In [8], we analyzed a number of images for the four objects of interest (bags, trolleys, persons, and groups of people), and

propose an effective feature set capable for discriminating the four classes with a high detection rate and a low false alarm rate. The features in the set represent the main statistics of geometric primitives for an object such as: corners, lines, circles, and other related statistics [8].

We follow the same approach and extract these features with the addition of the fitting ellipsis aspect ratio and the dispersion of the object. The fitting ellipse aspect ratio is calculated as the ratio between the length of minor axes and the length of major axes of the fitting ellipse. We further calculate the perimeter (the length of the external contour) and the area (the area under the external contour). The *dispersion* of an object is calculated as the ratio between the square of the perimeter and the area of the object. A full list of the features is illustrated in Table 1 and further described in [8]. Moreover, Figure 1 (right column) shows such features as extracted in a number of images.

**Table 1.** List of statistics of geometric primitives' features.

Corners	Circles	Lines	Other features
- No. of corners.	- No. of circles.	- No. of lines (strong, intermediate, and weak).	- Bounding box dispersion & Height/Width ratio
- The ratios and percentages between corners.	- The ratios and percentages between circles.	- No. of horizontal, vertical, diagonal lines, and ratios between them.	- Fitting ellipse aspect ratio
- Horizontal and vertical StDev.	- Horizontal and vertical StDev.		- Object dispersion

### 3. CLASSIFICATION

The classifiers that have been used for the classification experiments in our system are the Bayesian classifier BayesNet, C4.5 or Decision Trees, Sequential Minimal Optimization (SMO) algorithm [13], and MultiBoostAB (a variant of AdaBoost combining wagging and boosting) [12]. The performance of the classifier is evaluated in terms of classification accuracy (or detection rate for each class) and false positive rate (FPR). Classification accuracy is calculated as the proportion of the number of objects correctly detected against the total number of objects. The false positive rate is calculated as the proportion false positives against the sum of true negatives and false positives.

### 4. EXPERIMENTAL RESULTS AND ANALYSIS

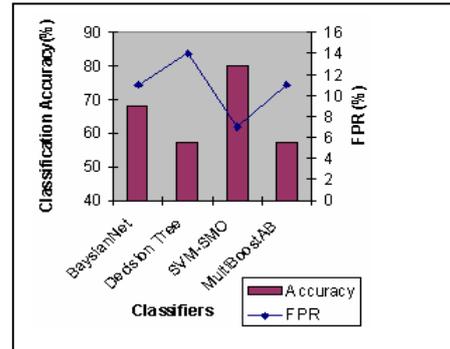
Experiments are conducted in order to compare different feature extraction techniques and evaluate their performance across a number of classifiers. For this purpose, we collected 600 images of trolleys, bags, single persons, and groups of people. These images were collected from video footage provided by our industrial partner and were taken in a number of airports around the world. Objects of interest in these images appear from different viewpoints, under different illumination conditions and in varying size and scale. We divided the images into two data sets: training set (400 images) and testing set (200 images), with equal number of images for each class.

For approach 1 and approach 2, we first extract SIFT keypoints and then apply PCA in order to reduce the dimensionality. In approach 1, we apply the majority rule described in Section 2.1.2 and then feed the results to the four different classifiers mentioned in previous section. For approach 2, a histogram is built for the reduced dimensions and the results are also fed to the multiple classifiers. Finally, for our approach (approach 3), we extract lines, circles, corners, and all other related statistical features and also feed them to the same classifiers. The results of classification based on these approaches are presented in Table 2, where classification accuracy and FPR are presented as a range across multiple classifiers, from the minimum to the maximum percentages.

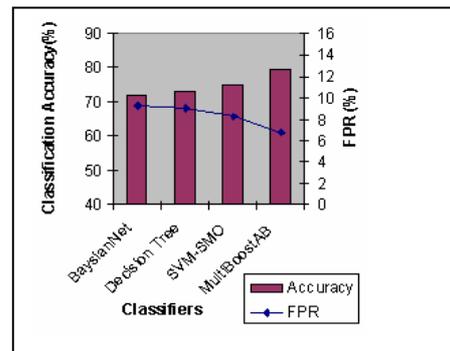
**Table 2.** Classification results as a range across multiple classifiers for the three approaches.

	Classification Accuracy	False Positive Rate
1 - SIFT keyp.	38% - 44.5%	20.6% - 22.8%
2 - SIFT hist [1]	57.5% - 80%	7% - 14%
3 - Our approach	72% - 79.5%	6.8% - 9.3%

It is clear from Table 2 that building a histogram for the SIFT keypoints outperforms the majority rule approach. The integral and non-local nature of the histogram as a feature results in a higher performance.



(a) Approach 2: SIFT keypoint histograms.



(b) Approach 3: Statistics of geometric primitives.

**Figure 2(a-b).** Classification accuracy and false positive rates for different approaches across a number of classifiers.

Moreover, by looking at Table 2, we also observe that the highest performance is achieved by our approach (approach 3), which is based on statistics of geometric primitives. This can be explained with the fact that in wide-area video surveillance, objects are often limited in size, and most often are low in texture and appear under different viewpoints. This results in a low number of detected SIFT keypoints and inconsistency of these keypoints across each class, leading to a lower classification performance compared to a classifier that is based on statistics of geometric primitives features.

In Figure 2, we plot the performance (classification accuracy and FPR) of the best two approaches (approach 2 and approach 3) across different classifiers. It is clear how the performance achieved based on geometric primitives' features proves better across a range of classification algorithms compared to the second-best approach (72%±79.5% accuracy vs. 57.5%±80%; alongside an FPR of 6.8%±9.3% vs. 7%±14%).

**Table 3.** The classification results of our approach with different datasets

	Average classification Accuracy	Average FPR
Original Dataset	74.86 %	8.4 %
Mixed Dataset	73.02 %	9 %

We have also experimented the invariance of our approach to different data sets. We experimented with the same four classifiers using a *mixed dataset* that includes the original 600 images (*original dataset*) with the inclusion of 124 images for the objects of interest that were collected from WWW. The mixed dataset is then divided into training dataset and testing dataset (2/3 (training) and 1/3 (testing)). The results of average classification accuracy and average FPR across the classifiers for our approach with the two datasets is presented in Table 3. From Table 3, we conclude that the performance of our approach is stable under different datasets. The results also imply that we can apply our approach to various environments and conditions without the need to re-tailor it.

## 5. CONCLUSION

In this paper, we compared three different approaches for classification that use different techniques for feature extraction. Based on the experimental results we obtained, we conclude that (i) the results of our approach for classification based on statistics of geometric primitives outperforms the other two approaches that are based on SIFT keypoints, (ii) the performance achieved by our approach is more invariant to the different classification learning methods compared to the other approaches and (iii) the performance of our approach is also stable under different datasets.

The results of our approach are encouraging considering the challenges inherent to the intra-class shape variation, illumination changes, variable viewpoints, and clutter. We plan in the future to experiment with other feature reduction methods, possibly Kernel Principle Component Analysis (KPCA), to improve the classification performance even further.

## 6. ACKNOWLEDGMENTS

This research is supported by the Australian Research Council and iOmniscient Pty Ltd under the ARC Linkage Project Grant Scheme 2006 - LP0668325

## 7. REFERENCES

- [1] B. Ayres and M. Boutell, "Home interior classification using SIFT keypoint histograms", *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, Minnesota, USA, pp. 1-6, 2007.
- [2] M.D. Beynon et al., "Detecting abandoned packages in a multi-camera video surveillance system", *Proc. of the IEEE Advanced Video and Signal Based Surveillance Conference*, Florida, USA, pp. 221-228, 2003.
- [3] L.M. Brown, "View independent vehicle/person classification", *Proc. of ACM 2nd International Workshop on Video surveillance & Sensor Networks*, New York, USA, pp. 114-123, 2004.
- [4] R. Fergus, P. Perona, A. Zisserman, "Object class recognition by unsupervised scale-invariant learning", *Proc. of IEEE Conference on Computer Vision and Pattern Recognition*, Wisconsin, USA, pp. 264-271, 2003.
- [5] T. Kadir, A. Zisserman, and M. Brady, "An affine invariant salient region detector", *Proceedings of European Conference on Computer Vision*, Prague, Czech Republic, pp. 228-241, 2004.
- [6] D.G. Lowe, "Distinctive image features from scale-invariant keypoints", *International Journal of Computer Vision*, Vol. 60, No. 2, pp. 91-110, 2004.
- [7] K. Mikolajczyk and C. Schmid, "A performance evaluation of local descriptors", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Vol. 27, No. 10, pp. 1615-1630, 2005.
- [8] A.F. Otoom, H. Gunes, and M. Piccardi, "Towards automatic abandoned object classification in visual surveillance systems", *Proceedings of Asia-Pacific Workshop on Visual Information Processing*, Tainan, Taiwan, pp. 143-149, 2007.
- [9] L. Sijun et al., "Detecting unattended packages through human activity recognition and object association", *Pattern Recognition*, Vol. 40, No. 8, pp. 2173-2184, 2007.
- [10] M. Tsuchiya, and H. Fujiyoshi, "Evaluating feature importance for object classification in visual surveillance", *Proc. of IEEE International Conference on Pattern Recognition*, Hong Kong, pp. 978-981, 2006.
- [11] G. Wang, Y. Zhang, and L. Fei-Fei, "Using dependent regions for object categorization in a generative framework", *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, New York, USA, pp. 1597-1604, 2006.
- [12] G.I. Webb, "MultiBoosting: a technique for combining boosting and wagging", *Machine Learning*, Vol. 40, No. 2, pp. 159-196, 2000.
- [13] H. Witten and E. Frank, *Data mining: Practical machine learning tools with java implementations*, Morgan Kaufmann, San Francisco, CA, 2000.
- [14] Non-Motion Detection Technology: [http://www.iomniscient.com/products\\_iq140.htm](http://www.iomniscient.com/products_iq140.htm) (Access date: 14 January 2008).

© [2008] IEEE. Reprinted, with permission, from [ Ahmed Fawzi Otoom, Hatice Gunes, Massimo Piccard, FEATURE EXTRACTION TECHNIQUES FOR ABANDONED OBJECT CLASSIFICATION IN VIDEO SURVEILLANCE , 15th IEEE International Conference on Image Processing, 2008. ICIP 2008]. This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the University of Technology, Sydney's products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org). By choosing to view this document, you agree to all provisions of the copyright laws protecting it