

# Detection of regular objects in baggage using multiple X-ray views

D Mery, G Mondragon, V Riffo and I Zuccar

Submitted 24.08.12

Accepted 05.11.12

*In order to reduce the security risk of a commercial aircraft, passengers are not allowed to take certain items in carry-on baggage. For this reason, human operators are trained to detect prohibited items using a manually controlled baggage screening process. In this paper, we propose the use of a method based on multiple X-ray views to detect some regular prohibited items with very defined shapes and sizes. The method consists of two steps: 'structure estimation', to obtain a geometric model of the multiple views from the object to be inspected (baggage); and 'parts detection', to detect the parts of interest (prohibited items). The geometric model is estimated using a structure from a motion algorithm. The detection of the parts of interest is performed by an ad-hoc segmentation algorithm (object dependent) followed by a general tracking algorithm based on geometric and appearance constraints. In order to illustrate the effectiveness of the proposed method, experimental results on detecting regular objects – razor blades and guns – are shown yielding promising results.*

Keywords: X-ray testing, baggage screening, luggage scan, multiple view imaging, computer vision.

## 1. Introduction

Since 9/11, aviation security screening with X-ray scanners has become a very important issue in airports. The inspection process, however, is complex because threatening items are very difficult to detect when placed in close-packed bags, superimposed by other objects and/or rotated showing an unrecognisable view<sup>[1]</sup>. In baggage screening, where human security plays an important role and inspection complexity is very high, human inspectors are still used. Nevertheless, during rush hours in airports, human screeners have only a few seconds to decide whether a bag contains a prohibited item or not, and detection performance is only about 80-90%<sup>[2]</sup>.

For these reasons, digital imaging and computer vision techniques have been developed in order to increase the effectiveness and automation of the inspection task. Before 9/11, however, the X-ray analysis of luggage mainly focused on capturing the images of their contents: the reader can find in<sup>[3]</sup> an interesting analysis carried out in 1989 of several aircraft attacks in the world, and the existing technologies to detect the terrorist threats based on thermal-neutron activation (TNA), fast-neutron activation (FNA) and dual-energy X-rays (used in medicine since early 1970). In the 1990s, explosive detection systems (EDS) were developed based

on X-ray imaging<sup>[4]</sup> and computed tomography through elastic scatter X-ray (comparing the structure of irradiated material against stored reference spectra for explosives and drugs)<sup>[5]</sup>.

All these works were concentrated on image acquisition and simple image processing, but they lack advanced image analysis to improve the detection performance. Nevertheless, the 9/11 attacks increased the security policies at airports, which also produced an interest in the scientific community for researching topics related to security using advanced computational techniques. In the last decade, the main contributions were: analysis of human inspection<sup>[6]</sup>, pseudo-colouring of X-ray images<sup>[7]</sup>, enhancement and segmentation of X-ray images<sup>[8]</sup> and detection of threatening items in X-ray images based on texture features (detecting a 9 mm Colt Beretta machine pistol)<sup>[9]</sup>, neural networks and fuzzy rules (yielding about 80% of performance)<sup>[10]</sup> and an SVM classifier (detecting guns in real time)<sup>[11]</sup>.

Recently, some algorithms based on multiple X-ray views were reported in the literature. For example: synthesis of new X-ray images obtained from kinetic depth effect X-ray (KDEX) images based on SIFT features in order to increase the detection performance<sup>[12]</sup>; active vision with X-ray, which allows modifying the viewpoint of the target object in order to obtain better X-ray images to analyse (detecting razor blades in different cases)<sup>[13]</sup>; and tracking across multiple X-ray views in order to verify the diagnoses performed using a single view<sup>[14]</sup>. The key idea of this method is: (i) to segment potential parts (regions) of interest in each view using an application-dependent method that analyses 2D features in each single view, ensuring the detection of the object parts of interest (not necessarily in all views) and allowing false detections; (ii) to match and track the potential regions based on similarity and geometrical multiple view constraints, eliminating those that cannot be tracked; and (iii) to analyse the tracked regions, including those views where the segmentation fails (the positions can be predicted by re-projection). This algorithm will be explained in the next section in further detail because it is the core of this paper.

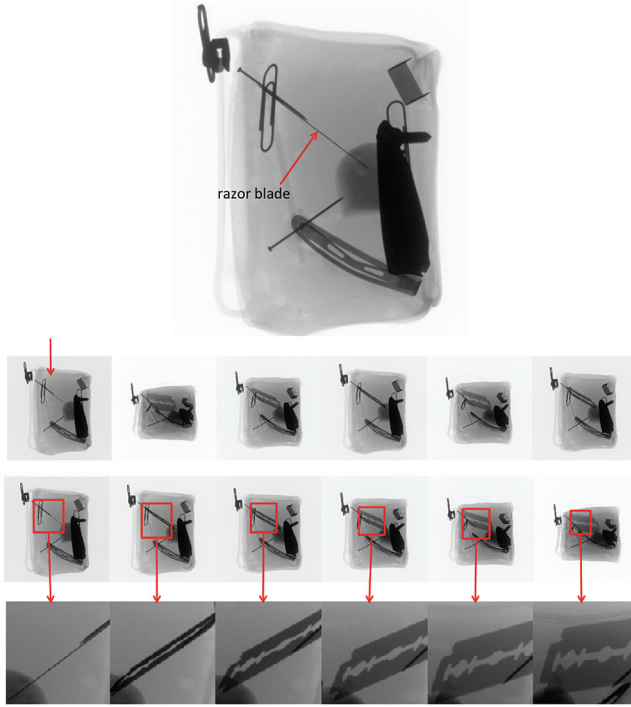
In baggage screening, the use of multiple view information yields a significant improvement in performance because certain items are difficult to recognise using only one viewpoint, as we illustrate in Figure 1, where we detected a razor blade in a pencil case using the proposed method. It is clear that this detection performance could not be achieved with only the first view of the sequence.

In this work, we use the general methodology proposed by us in<sup>[14]</sup> and implement algorithms to automatically detect regular objects in baggage (such as razor blades and guns) with multiple X-ray views. We show the robustness of the approach against poor segmentation or noise because these false detections are not attached to the object and therefore they cannot be tracked.

The rest of the paper is organised as follows: the multiple view approach is summarised in Section 2; the *ad-hoc* single view detectors are explained in Section 3; the results obtained in several experiments are shown in Section 4; and some concluding remarks are given in Section 5.

Domingo Mery, German Mondragon and Irene Zuccar are with DCC – Pontificia Universidad Católica de Chile. Email: dmery@ing.puc.cl / german.mondragon@gmail.com / irene.zuccar@usach.cl

Vladimir Riffo is with DCC – Pontificia Universidad Católica de Chile and DIICC – Universidad de Atacama. Email: vriffo1@uc.cl



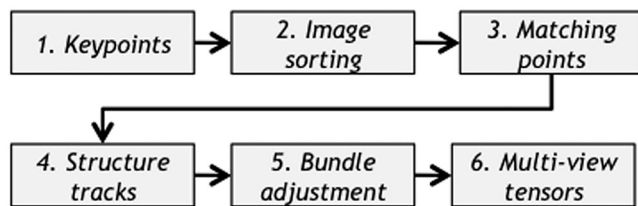
**Figure 1. Detection of a razor blade in a pencil case using our approach: First X-ray image of the sequence, 1430 × 900 pixels (top); Unsorted and sorted sequences with six images (middle); Detection in each image (bottom)**

## 2. Multiple view approach

In this section we summarise the multiple view approach outlined in<sup>[14]</sup> using *ad-hoc* single view detectors for regular objects. The proposed method follows two main steps: ‘structure estimation’, to obtain a geometric model of the multiple views from the object itself; and ‘parts detection’, to detect the object parts of interest.

### 2.1 Structure estimation

The approach outlined in this section is based on well-known structure from motion<sup>[15]</sup> estimated from a sequence of  $m$  images taken from a rigid object at different viewpoints (see Figure 2).



**Figure 2. Block diagram of structure estimation**

The original image sequence is stored in  $m$  images  $J_1, \dots, J_m$ . For each image, SIFT keypoints are extracted<sup>[16]</sup>. Thus, not only a set of 2D image positions,  $x$ , but also descriptors,  $y$ , are obtained. The images of the sequence are sorted using a visual vocabulary tree in order to obtain a sequence with small changes between consecutive frames<sup>[17]</sup>, as shown in Figure 1. For two consecutive and sorted images,  $I_i$  and  $I_{i+1}$ , SIFT keypoints are matched using the algorithm suggested by Lowe<sup>[16]</sup>, which rejects matches that are too ambiguous. Afterwards, the fundamental matrix between views  $i$  and  $i + 1$ ,  $F_{i,i+1}$ , is estimated using RANSAC<sup>[15]</sup> to remove outliers. We look for all possible structure tracks with one keypoint in each

image of the sequence.

The determined tracks define  $n$  image point correspondences over  $m$  views. They are arranged as  $x_{i,j}$  for  $i = 1, \dots, m$  and  $j = 1, \dots, n$ . Bundle adjustment estimates 3D points  $\hat{X}_j$  and camera matrices  $P_i$  so that  $\sum \|x_{i,j} - \hat{x}_{i,j}\|^2$  is minimised, where  $\hat{x}_{i,j}$  is the projection of  $\hat{X}_j$  by  $P_i$ . A RANSAC approach is used to remove outliers. Bundle adjustment<sup>[15]</sup> provides a method for computing bifocal and trifocal tensors from projection matrices  $P_i$ , which will be used in the next section.

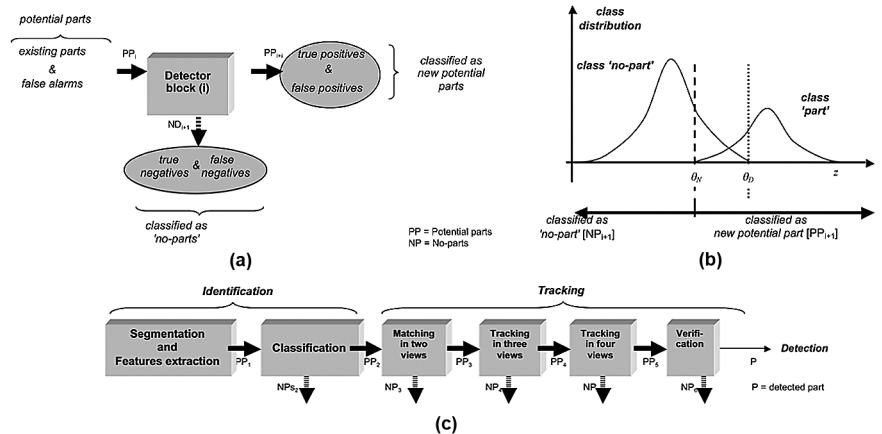
### 2.2 Parts detection

In this section we give details of the algorithm that detects the object parts of interest. The algorithm consists of the following two main steps: identification and tracking, as shown in Figure 3.

The strategy is to ensure the detection of the ‘existing parts of interest’ in the first step, allowing the inclusion of ‘false alarms’. The discrimination between both is achieved in the second step using multiple view analysis, where an attempt is made to track the potential parts of interest along the image sequence.

In the identification, potential parts of interest are segmented and classified in each image  $I_i$  of the sequence. It is an *ad-hoc* single view detector that depends on the application. In Section 3, two algorithms will be explained for the detection of regular objects (razor blades and guns).

An existing part of interest can be successfully tracked in the image sequence because its appearance in the images is similar and their projections are located in the positions dictated by geometric conditions. In contrast, false alarms can be successfully eliminated



**Figure 3. Block diagram of parts detection: each block separates the new potential parts from the no-parts (a) according to the class distribution (b). The whole diagram follows a cascade schema as shown in (c)**

in this manner, since they do not appear in the predicted places on the following images and, thus, cannot be tracked. The tracking in the image sequence is performed using algebraic multi-focal constraints: bifocal (epipolar) and trifocal constraints among others<sup>[15]</sup> obtained from projection matrices estimated in the previous step outlined in Section 2.1, where the geometric model is obtained from the target object itself.

Each sub-step  $i$  of the automated multiple view analysis can be understood as a detector block, as shown in Figure 3. The potential parts ( $PP_i$ ) consisting of existing parts and false alarms are classified as either new potential parts ( $PP_{i+1}$ ) or no-parts ( $NP_{i+1}$ ) (Figure 3(a)). In a training phase, each detector block is tuned so that the maximal number of false alarms is eliminated from the potential parts without discriminating the existing defects (see  $\theta_s$  in Figure 3(b)). The throughput cycle can be considerably incremented if we use an additional decision boundary (see  $\theta_d$  in Figure 3(b)), which guarantees the detection of defects in previous stages without computing the next steps.

The reader is referred to<sup>[14]</sup> for a detailed description of the tracking algorithm.

### 3. Object-dependent single view detector

An object-dependent algorithm must be defined to detect automatically potential parts of interest in a single test image. As mentioned above, in order to test our method, we developed two algorithms that are able to detect razor blades and guns. In this section, they will be explained in further detail.

#### 3.1 Detection of razor blades

The algorithm is based on the matching of SIFT keypoints<sup>[16]</sup> and was proposed by us in the first part of<sup>[13]</sup> for active vision. In our approach, we use a SIFT description of the target object in all feasible poses by rotating two axes in nine steps. All extracted descriptors are stored in an arrangement  $\mathbf{P}$ , where  $\mathbf{p}_j$  means the  $j$ -th descriptor, for  $j = 1, \dots, m$ . Each descriptor  $\mathbf{p}_j$  has a corresponding pose  $r_j$ . In our example,  $r_j \in [1, 81]$  for  $9 \times 9$  poses. All SIFT descriptors of the test image of the inspection object are extracted and stored in an arrangement  $\mathbf{Q}$ , where  $\mathbf{q}_i$  means the  $i$ -th descriptor of the test image for  $i = 1, \dots, n$ . Now, all duplets  $(\mathbf{q}_i, \mathbf{p}_j)$  that fulfill the condition  $\|\mathbf{q}_i - \mathbf{p}_j\| < \theta_e$  for  $i = 1, \dots, n$  and  $j = 1, \dots, m$  are selected, where  $\theta_e$  is a minimum distance threshold, and  $\|\mathbf{q}_i - \mathbf{p}_j\|$  means the Euclidean distance between both vectors. Afterwards, for each selected descriptor the corresponding pose  $r_j$  is obtained. The selected descriptors and their corresponding poses will be stored in  $\mathbf{Q}$  and  $\mathbf{R}$ , respectively. Thus, we have (i)  $\mathbf{Q}$ : all keypoints of the test image that have been matched with keypoints of the target object, and (ii)  $\mathbf{R}$ : the corresponding poses for the selected keypoints  $\mathbf{Q}$ .

The detection is performed in the following two steps:

- (i) Clustering: in  $\mathbf{Q}$  we find all keypoints of the same pose that are close to each other in the test image. Thus, we define subwindows  $\mathbf{W}_B$  that have at least  $\theta_B$  keypoints of the same pose. In our experiments, we set the size of  $\mathbf{W}_B$  equal to  $80 \times 80$  pixels, and  $\theta_B = 3$ .
- (ii) Merging: all subwindows  $\mathbf{W}_B$  that are connected or overlapped will be merged in a new larger subwindow  $\mathbf{W}_G$ . The subwindow that encloses the highest number of keypoints of the same pose will be selected if this number is equal to or greater than  $\theta_G$ , ensuring at least  $\theta_G$  descriptors of the same pose in the selected window. In our experiments, we set  $\theta_G = 2$  in order to ensure the detection (allowing false alarms). The selected subwindow will be called  $\mathbf{W}_S$  and will correspond to a potential target object. If no subwindow fulfills this condition, then no potential target object is detected.

#### 3.2 Detection of guns

In the computer vision community, many object detection and classification problems have been recently solved – without segmentation – using sliding windows. Sliding window approaches have established themselves as the state-of-the-art in computer vision problems, where an object must be separated from the background (see, for example, successful applications in face detection<sup>[18]</sup>). In sliding window methodology, a detection window is passed over an input image in both horizontal and vertical directions and, for each localisation of the detection window, a classifier decides to which class the corresponding portion of the image belongs, according to its features. Multiple detection can be eliminated using non-maximum suppression<sup>[18]</sup>.

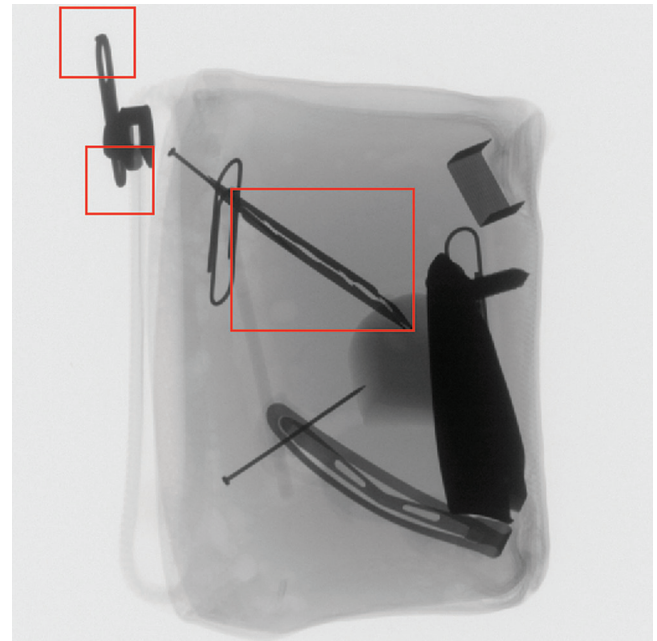
We used sliding windows to detect guns; however, since there are many types of gun, our approach to detect a gun is based on the detection of its trigger. We observed that the shape of triggers has a smaller variability in comparison to the shape of the guns. For this reason we collected 100 images of guns from Google Images and cropped their triggers. Afterwards, we built a dataset with positive classes (trigger images) and negative classes (no trigger images). We trained a classifier using this dataset. We tested with several features and classifiers. A simple and fast solution was achieved using a Mahalanobis distance classifier with seven geometric features (a Hu moment, a Fourier descriptor, centre of mass, and minor and major axes of a fitted ellipse).

### 4. Experimental results

We experimented on X-ray images from two different objects: (i) detection of razor blades; and (ii) detection of guns.

#### 4.1 Experiments on razor blades

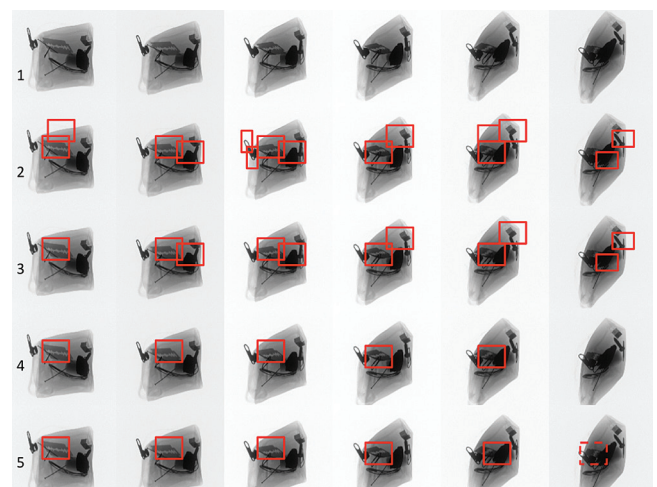
We tested four sequences of razor blades with four-six images with very good results. Figure 4 illustrates the detection using the single view detector outlined in Section 3.1. We observe that the razor blade was identified; however, there are two false alarms. They were filtered out after tracking steps.



**Figure 4. Identification of potential razor blades in a single view: there are two false alarms (upper left) and a true positive (large rectangle in the middle)**

Figure 5 shows the results obtained in each step. We can see that the razor blade was not identified by the single view detector in the last image; however, after tracking using the geometric model, it was possible to re-project its position in this view (see the last dashed rectangle). Another example is illustrated in Figure 1; again, the razor blade was not identified in the first image, however, it could be re-projected.

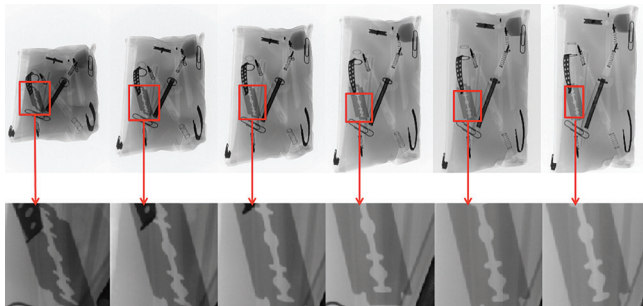
Another experiment with the same result is shown in Figure 6. In this experiment there were 35 potential razor blades identified in



**Figure 5. Results in each step: (1) sorted image sequence; (2) detection of potential razor blades using the single view detector; (3) remaining potential razor blades after matching in two views; (4) after three views; (5) four views**



the six-image sequence (only six of them were real existing razor blades); however, after using the tracking algorithm, the 29 false alarms were eliminated.



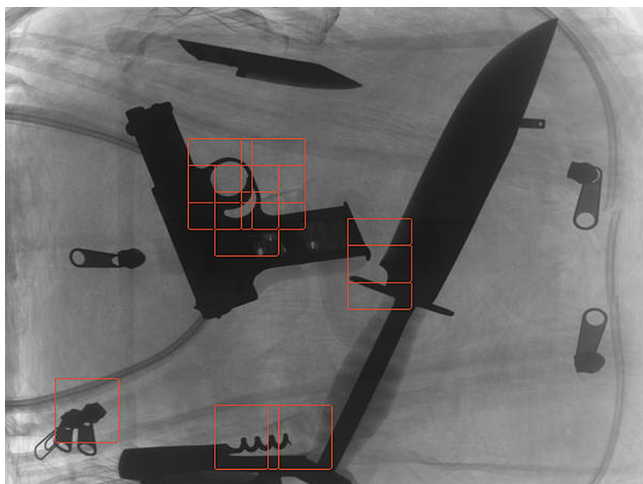
**Figure 6.** Detection of a razor blade in a pencil case. Top: sequence with six X-ray images, 1430 × 900 pixels; Bottom: detection

#### 4.2 Experiments on guns

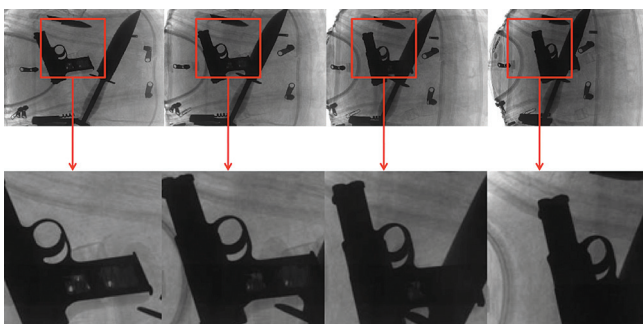
We tested ten sequences of four-five images of bags and backpacks containing a gun. The detection was achieved in sequences where the trigger was distinguishable. An example of the single view detector is shown in Figure 7. In this case the multiple false alarms were filtered out by the tracking algorithm, as shown in Figure 8. Another experiment that shows a very good detection with some occlusion is illustrated in Figure 10. Nevertheless, in intricate sequences (see, for example, Figure 9) the gun could not be detected because it was too occluded.

#### 4.3 Performance

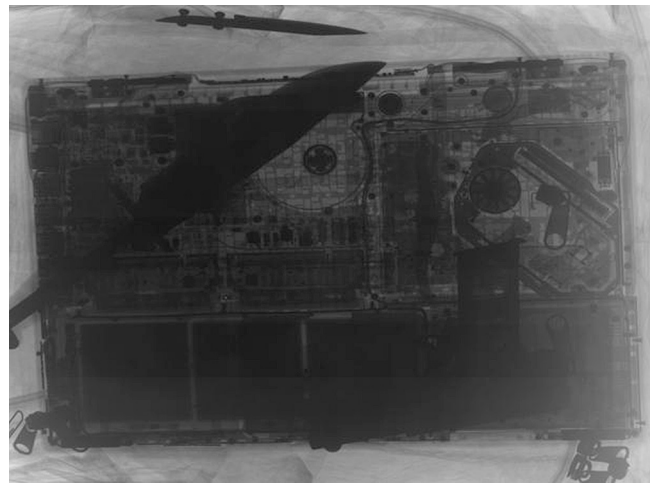
Tables 1 and 2 summarise the results on razor blades and guns with 14 sequences (64 X-ray images). Some of them are illustrated in the mentioned Figures.  $m$  corresponds to the number of images in



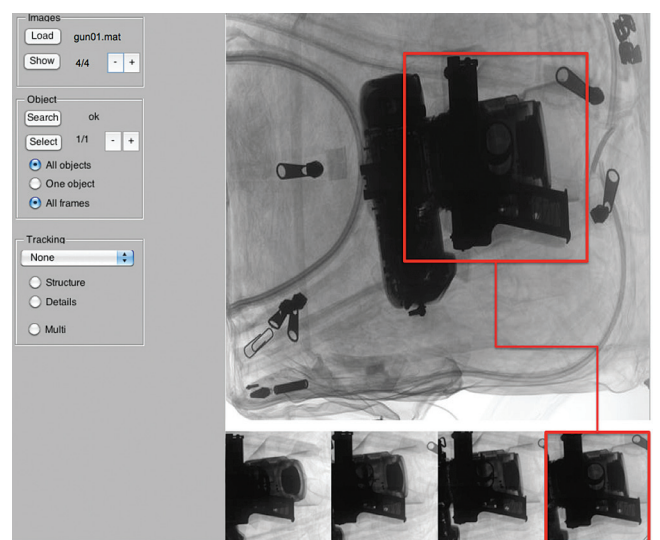
**Figure 7.** Single view detection of a gun: we observe that there are several false alarms that will be eliminated after tracking, as shown in Figure 8



**Figure 8.** Detection of a gun in a bag. Top: sequence with four X-ray images, 452 × 612 pixels; Bottom: detection



**Figure 9.** X-ray image of a gun and a knife on a laptop: the detection of the gun was impossible. See performance statistics in Table 2, Sequence 9



**Figure 10.** Developed graphic user interface (GUI) showing the detection of an occluded gun

the sequence.  $\text{SIFT}/m$  means the average of the number of SIFT keypoints extracted per image. BA is the number of structure tracks found by the bundle adjustment algorithm.  $n_1$  is the number of segmented potential regions in the whole image sequence, and  $n_1/m$  is the average of segmented regions per image.  $n_l$  is the number of  $l$ -tuples tracked in the sequence.  $n_d$  is the number of detected parts. GT is the number of existing parts (ground truth). FP and TP are the number of false and true positives after eliminating multiple overlapped detections. Ideally,  $\text{FP} = 0$  and  $\text{TP} = \text{GT}$ . If we include all sequences, the average performance is given by: Precision =  $\text{TP}/(\text{TP} + \text{FP}) = 70\%$  and Recall =  $\text{TP}/\text{GT} = 86\%$ . Nevertheless, if we exclude the last two gun sequences (they are not allowed in baggage screening because laptops must be removed from bags) Precision = 86% and Recall = 100%.

#### 4.4 Implementation

We implemented our approach in a Matlab graphic user interface (Figure 10). We used the implementation of SIFT, visual vocabulary, etc from VLFeat<sup>[19]</sup>. The rest of algorithms were implemented in Matlab. For multiple view matching,  $\epsilon_2 = 30$  pixels and  $\epsilon_3 = 50$  pixels. The computing time depends on the application, however, in order to present a reference, for Figure 1 the results were obtained after 30 s on an iMac OS X 10.6.6, processor 3.06 GHz Intel Core 2 Duo, 4 GB RAM memory. The code of the Matlab implementation is available on our webpage<sup>[20]</sup>.

**Table 1. Detection of razor blades\***

Seq	Size	$m$	SIFT/ $m$	BA	$n_1$	$n_1/m$	$n_2$	$n_3$	$n_4$	$n_q$	$n_d$	GT	FP	TP
1	1430 × 900	6	2372	30	35	6	18	4	1	1	1	1	0	1
2	850 × 850	6	1679	4	14	2	13	8	1	1	1	1	0	1
3	850 × 850	6	1312	2	12	2	12	4	1	1	1	1	0	1
4	1430 × 900	4	5135	58	26	7	15	6	2	2	2	1	1	1
Total	–	22	–	–	–	17	–	–	–	–	5	4	1	4

\* Variables used in this Table are explained in Section 4.3

**Table 2. Detection of guns\*\***

Seq	Size	$m$	SIFT/ $m$	BA	$n_1$	$n_1/m$	$n_2$	$n_3$	$n_4$	$n_q$	$n_d$	GT	FP	TP
1	459 × 612	4	2226	24	73	18	268	162	66	5	5	1	0	1
2	459 × 612	5	2253	6	114	23	573	347	164	15	15	1	0	1
3	459 × 612	4	2222	39	44	11	171	71	32	5	5	1	0	1
4	459 × 612	4	2242	35	38	10	162	87	8	2	2	1	1	1
5	459 × 612	4	2192	113	33	8	182	192	91	8	8	1	0	1
6	459 × 612	4	2297	39	88	22	596	166	48	7	7	1	0	1
7	459 × 612	4	662	33	103	26	1058	1407	1108	38	38	1	1	1
8	459 × 612	4	662	33	8	2	14	9	3	1	1	1	0	1
9	459 × 612	5	2246	162	180	36	3041	2916	1221	71	71	1	2	0
10	459 × 612	4	1473	62	93	23	600	509	376	17	17	1	1	0
Total	–	42	–	–	–	179	–	–	–	–	169	10	4	8

\*\* Variables used in this Table are explained in Section 4.3

## 5. Conclusions

In this paper we presented the use of a generic methodology that can be used to detect regular prohibited items (such as razor blades and guns) in baggage automatically, yielding promising results. The proposed approach is an application of state-of-the-art computer vision techniques. It filters out false positives resulting from segmentation steps performed on single views of an object by corroborating information across multiple views.

Using multiple views (instead of one) the matching accuracy and robustness (*ie* tolerance to false-positive detections) of the detection of physical features on an object is increased. The detection method is image-based (2D appearance-based detection). By using multiple views, the method is able to increase the detection rates and robustness of 2D feature detection, in comparison to application of the same method in a single image. We believe that our methodology is a useful alternative for assisting human operators in baggage screening.

## Acknowledgements

This work was supported by grant Fondecyt No 1100830 from CONICYT – Chile.

## References

1. G Zentai, 'X-ray imaging for homeland security', IEEE International Workshop on Imaging Systems and Techniques (IST 2008), pp 1-6, September 2008.
2. S Michel, S Koller, J de Ruiter, R Moerland, M Hogervorst and A Schwaninger, 'Computer-based training increases efficiency in X-ray image interpretation by aviation security screeners', Security Technology 2007, 41st Annual IEEE International Carnahan Conference, pp 201-206, October 2007.
3. E Murphy, 'A rising war on terrorists', Spectrum, IEEE, Vol 26, No 11, pp 33-36, November 1989.
4. N Murray and K Riordan, 'Evaluation of automatic explosive detection systems', Proceedings of the 29th Annual International Carnahan Conference on Security Technology 1995, Institute of Electrical and Electronics Engineers, pp 175-179, October 1995.
5. H Strecker, 'Automatic detection of explosives in airline baggage using elastic X-ray scatter', Medicamundi, Vol 42, pp 30-33, July 1998.
6. A Wales, T Halbherr and A Schwaninger, 'Using speed measures to predict performance in X-ray luggage screening tasks', 43rd Annual International Carnahan Conference on Security Technology 2009, pp 212-215, October 2009.
7. J Chan, P Evans and X Wang, 'Enhanced colour coding scheme for kinetic depth effect X-ray (KDEX) imaging', 2010 IEEE International Carnahan Conference on Security Technology (ICCST), pp 155-160, October 2010.
8. M Singh and S Singh, 'Optimising image enhancement for screening luggage at airports', Proceedings of the 2005 IEEE International Conference on Computational Intelligence for Homeland Security and Personal Safety, CIHSPS 2005, pp 131-136, 31 May - 1 April 2005.
9. C Oertel and P Bock, 'Identification of objects-of-interest in X-ray images', Applied Imagery and Pattern Recognition Workshop 2006, AIPR 2006, 35th IEEE, p 17, October 2006.
10. D Liu and Z Wang, 'A united classification system of X-ray image based on fuzzy rule and neural networks', 3rd International Conference on Intelligent System and Knowledge Engineering 2008, ISKE 2008, Vol 1, pp 717-722, November 2008.
11. S Nercessian, K Panetta and S Agaian, 'Automatic detection of potential threat objects in X-ray luggage scan images', 2008 IEEE Conference on Technologies for Homeland Security, pp 504-509, May 2008.
12. O Abusaeeda, J Evans, D Downes and J Chan, 'View synthesis of KDEX imagery for 3D security X-ray imaging', Proceedings of the 4th International Conference on Imaging for Crime Detection and Prevention (ICDP-2011), 2011.
13. V Rizzo and D Mery, 'Active X-ray testing of complex objects',

Insight, Vol 54, No 1, pp 28-35, 2012.

14. D Mery, 'Automated detection in complex objects using a tracking algorithm in multiple X-ray views', Proceedings of the 8th IEEE Workshop on Object Tracking and Classification Beyond the Visible Spectrum (OTCBVS 2011), in conjunction with CVPR 2011, Colorado Springs, pp 41-48, 2011.
15. R I Hartley and A Zisserman, Multiple View Geometry in Computer Vision, 2nd edition, Cambridge University Press, 2003.
16. D Lowe, 'Distinctive image features from scale-invariant keypoints', International Journal of Computer Vision, Vol 60, No 2, pp 91-110, 2004.
17. J Sivic and A Zisserman, 'Efficient visual search of videos cast as text retrieval', IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol 31, No 4, pp 591-605, 2009.
18. P Viola and M Jones, 'Robust real-time object detection', International Journal of Computer Vision, Vol 57, No 2, pp 137-154, 2004.
19. A Vedaldi and B Fulkerson, 'VLFeat: An open and portable library of computer vision algorithms', <http://www.vlfeat.org/>, 2008.
20. D Mery, 'BALU: A Matlab toolbox for computer vision, pattern recognition and image processing', <http://dmery.ing.puc.cl/index.php/balu>, 2011.