

# Optical Engineering

[SPIDigitalLibrary.org/oe](http://SPIDigitalLibrary.org/oe)

## **Flying target detection and recognition by feature fusion**

Levente Kovács  
Andrea Kovács  
Ákos Utasi  
Tamás Szirányi



**SPIE**

# Flying target detection and recognition by feature fusion

Levente Kovács  
Andrea Kovács  
Ákos Utasi  
Tamás Szirányi

Hungarian Academy of Sciences  
Computer and Automation Research Institute  
Distributed Events Analysis Research Laboratory  
Kende u. 13-17, 1111 Budapest, Hungary  
E-mail: [levente.kovacs@sztaki.mta.hu](mailto:levente.kovacs@sztaki.mta.hu)

**Abstract.** We present a near-real-time visual-processing approach for automatic airborne target detection and classification. Detection is based on fast and robust background modeling and shape extraction, while recognition of target classes is based on shape and texture-fused querying on a-priori built real datasets. The presented approach can be used in defense and surveillance scenarios where passive detection capabilities are preferred (or required) over a secured area or protected zone. © 2012 Society of Photo-Optical Instrumentation Engineers (SPIE). [DOI: [10.1117/1.OE.51.11.117002](https://doi.org/10.1117/1.OE.51.11.117002)]

Subject terms: object segmentation; target detection; classification.

Paper 120791P received Jun. 1, 2012; revised manuscript received Sep. 5, 2012; accepted for publication Oct. 1, 2012; published online Nov. 2, 2012.

## 1 Introduction

Visual detection, recognition, classification and tracking of stationary or moving targets are among the most active research areas in computer vision and image processing fields. Applications built on results of these research areas are constantly sought to be deployed for both defensive and offensive scenarios, including civilian and military use. For civilian applications, wide area surveillance, crowd and traffic monitoring, and target tracking are the most important fields, while for military applications, troops and asset protection, region of interest surveillance, target detection, and tracking are probably the most important scenarios. Aiding such tasks by intelligent and automatic visual processing is important since such methods can support the detection, recognition, and alerting tasks of security personnel. Also, visual processing sensors/nodes can provide a means for passive detection (without requiring active signals), thus making them harder to detect and disarm in case of sensitive scenarios.

This paper presents a solution for one aspect of the above-described wide range of possibilities, focusing on automatic airborne target detection and classification. The presented approach can be used in defense and surveillance scenarios, where passive detection capabilities are preferred (or required) over a secured area or protected zone. The goals are to automatically detect and recognize the class of observed flying targets from varying angles, views, size, and environmental conditions, while running on commodity hardware.

Lu et al.<sup>1</sup> presented a small-ship target detection method, where point-like infrared images of small ships are processed to automatically detect ships on the sea level from a distance. Simple edge detection on a median filtered image is used to extract possible ship locations. In other works<sup>2</sup> small targets above a sea or sky background are extracted by infrared processing by using directional derivative operators and clustering. Deng et al.<sup>3</sup> present small target detection in infrared, based on self-information maps and locally adaptive background thresholding and region growing, producing robust detection results. While infrared processing can help the detection task (especially during the night), it is not suitable

for generic classification because of low resolution and less visual information.

Some target detection (without tracking and classification) methods<sup>4,5</sup> present object detection based on multiscale color and saliency information on single frames/images, detecting outlier regions based on local features as target candidates with good results, but not suitable for real-time video processing. Lia et al.<sup>5</sup> present a visual missile-like target detection approach based on image segmentation, motion analysis for target region selection, and a target boundary extraction step, validating on videos captured from three-dimensional (3-D) simulations. The *K*-means-based segmentation and the histogram-based target region extraction is not suitable for our purposes, since we have highly dynamic backgrounds with changing light, moving clouds, and vapor trails with free camera motions, and the grayscale histograms do not contain enough information for such scenarios.

Elsewhere, low flying targets are segmented<sup>6</sup> above the sea-sky line, by first locating the skyline, then using neighborhood averaging and directional Sobel operators to enhance the object boundaries. Bibby et al.<sup>7</sup> present a color-based tracking approach on a sea background using a stabilized camera on a moving platform. They use color and gradient-based mean shift tracking, without object detection or classification. Such an approach could be integrated with our proposed method for tracking purposes.

Wenga et al. present a flying target detection and tracking method<sup>8</sup> in infrared. Here, the goal is detection and tracking, without recognition/classification. Image complexity, number of objects and number of other large areas (e.g., cloud objects) are taken into consideration, and detection is performed depending on the weather condition (clouds or clear skies). Also, clouds are separated based on histogram analysis, assuming white cloud color. On the other hand, our approach does not use or depend on such information, any background clutter (clouds, vapor trails, smoke, independent of their color) get automatically discarded based on their nonrelevance as target candidates (based on their features), and the method is independent on the presence of such clutter. Other infrared-based approaches<sup>9</sup> also exist for target detection and tracking, although somewhat constrained since it requires static cameras, without classification capabilities. Wang et al.<sup>10</sup> present infrared target recognition on aerial

imagery by a multifeature method, sensitive to various geometrical shapes (circles, lines, etc.) of ground targets. Blasch et al.<sup>11</sup> present an approach for visual and infrared target tracking for day and night applicability, with the goal of keeping the target IDs in cluttered environments for robust long-term tracking that can be used for robust tracking after the objects are detected and categorized.

Noor et al.<sup>12</sup> present a model generation approach for object recognition, using multiple views of objects to build a model database. They start with scale invariant feature transform (SIFT)<sup>13</sup> descriptors for relevant corner point extraction, used to build a region-neighborhood graph that is used for object matching. In our case, we needed more robust interest point extraction because of the variances in backgrounds and viewing angles; thus, we present and use a more robust point extraction approach.

In our previous work, we introduced small target detection<sup>14</sup> and flying target tracking.<sup>15</sup> The current paper builds on these results and presents a more thorough investigation concentrating on shape and texture fused target recognition with extensive evaluation.

The novelties of the work presented in this paper are:

- Using an object detection and extraction method based on an enhanced interest point detection approach, which is robust against noise and clutter in the scene (e.g., clouds, vapor trails, other interference-like illumination changes), which is a new approach with respect to classical multilayer background/foreground modeling methods, shifting the complexity from the background modeling to a faster process of robust boundary point, and contour segment extraction.
- Presenting results in recognition of the extracted object classes, based on the combination of their shape and texture information, using a simple and fast approach based on indexing using BK-trees<sup>16</sup> and a turning function based metric combined with MPEG (Motion Picture Experts Group)-7 texture features.
- As opposed to several other methods that use pre-segmented datasets for training a classifier, no manual segmentation or annotation is done during the training and the detection/recognition process in our case. The dataset for building the initial shape index is gathered automatically by the same algorithm used during recognition, by running it on videos containing the target classes. The used video dataset contains real-life captures of flying targets on real background, without any simulation. The goal was to concentrate on real-life applicability and real environmental properties (varying illumination, cluttered background, multiple simultaneous targets, etc.).

Figure 1 shows the main algorithmic steps of the proposed method. The first step of the approach applies a novel boundary points detection technique, which is detailed in Sec. 2. In the second step, we create a Markovian foreground-background separation method for efficient object extraction, which is described in Sec. 2.2. Finally, for target classification we extract shape and texture features from the segmented objects, which are described in Secs. 2.3 and 3.

We will present results and evaluation of the object detection, extraction and recognition phases, showing the viability of the approach.

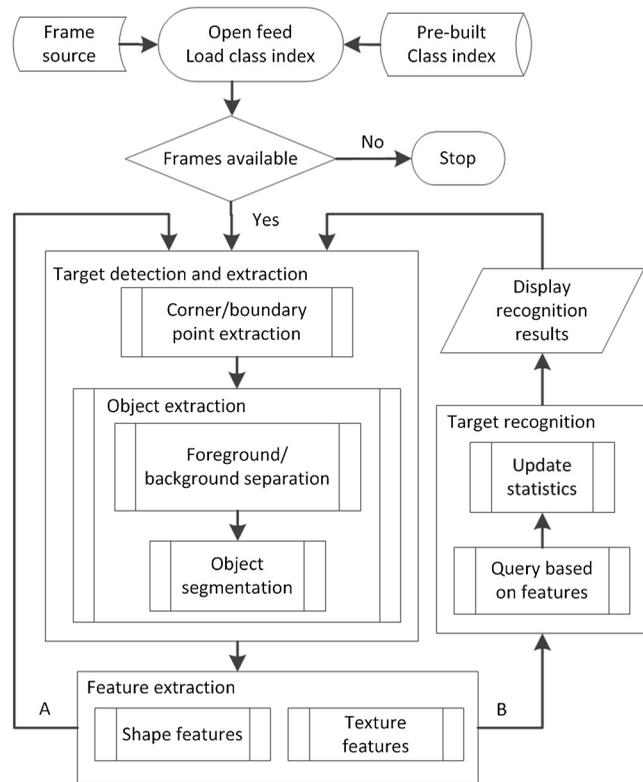


Fig. 1 Sequence diagram of the presented solution. Branches A and B run in parallel.

## 2 Target Detection and Feature Extraction

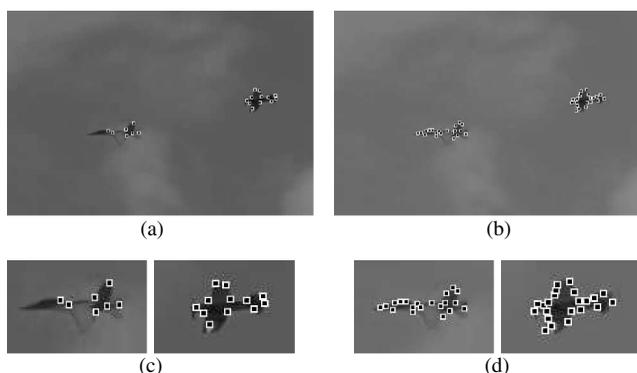
The first step towards the recognition phase is object extraction and object feature extraction. In this section, we will present these two steps. First, object extraction will be described based on a contour detection method that applies a novel interest point detector. Then, we present the object feature extraction step, used later on.

For the extraction of the object silhouettes, we propose a two-step probabilistic method, which achieves high computation performance, and works with multiple objects with different properties (e.g., size, color, shape, texture), on a changing cloudy sky background. Before introducing the proposed method, we list our assumptions:

- (1) Camera is not required to be static, but the consecutive video frames should contain overlapping parts.
- (2) Moving objects are smaller than the background part of the frames.
- (3) Background is not required to be completely homogeneous (e.g., clear skies), but should contain large homogeneous areas (which can be sky, clouds, vapor trails, etc.).

### 2.1 Feature Point Extraction

According to our assumptions, the aim is to find small (relative to the entire frame size) objects in front of a non-homogeneous background (that may contain clouds, sky regions, etc.) in the image  $I_t$  at time  $t$ . The first step of localizing foreground objects is to extract interest/feature points. The challenge in detecting such points is that the contours of the moving foreground objects are usually of low contrast,



**Fig. 2** Contour point detection: (a) original Harris corner detector;<sup>17</sup> (b) proposed modified Harris for edges and corners (MHEC) point detector; (c) and (d) respective objects zoomed.

and contain high curvature parts. Therefore, traditional point detectors, like the Harris corner detector,<sup>17</sup> cannot represent them accurately [see Fig. 2(a)]. To compensate for such drawbacks, we use a modification of the Harris detector, which was introduced earlier<sup>18</sup> and was applied for complex object contour recognition with parametric active contour algorithms.

The modified Harris for edges and corners (MHEC) method described here, adopts a modified characteristic function, and is able to emphasize not only corners, but edges as well, therefore it is more suitable for contour features than the traditional method.

The original Harris detector is based on the principle that intensity has large changes in multiple directions simultaneously at corner points. The method defines the  $R$  characteristic function for classifying image regions:

$$R = \text{Det}(M) - k * \text{Tr}^2(M), \quad (1)$$

where  $\text{Det}$  is the determinant and  $\text{Tr}$  is the trace of the  $M$  Harris matrix [see Eq. (2)], and  $k$  is a coefficient, usually chosen around 0.04. The  $M$  Harris matrix is defined as:

$$M = \begin{bmatrix} A & C \\ C & B \end{bmatrix}, \quad (2)$$

where  $A = \dot{x}^2 \otimes w$ ,  $B = \dot{y}^2 \otimes w$ , and  $C = \dot{x} \dot{y} \otimes w$ , with  $\dot{x}$  and  $\dot{y}$  denoting the approximation of the first order derivatives, and  $w$  is a Gaussian window.

$M$  describes the shape at an image point and its eigenvalues (denoted by  $\lambda_1$  and  $\lambda_2$ ) give a rotation invariant characterization of the curvatures in the small neighborhood of the point. Eigenvalues separate different regions; both of them are large in corner regions, only one of them is large in edge regions, and both of them are small in homogeneous (flat) regions.

In our case, when emphasizing edge and corner regions simultaneously, we exploit the fact that they both have one large eigenvalue, therefore  $L = \max(\lambda_1, \lambda_2)$  is able to separate homogeneous and nonhomogeneous regions in the image.<sup>18</sup> Let  $b_i = \{[x_i - r, x_i + r] \times [y_i - r, y_i + r]\}$  mark a window surrounding a specific  $p_i = (x_i, y_i)$  pixel.  $p_i$  is the element of the  $C$  contour point set, if it satisfies the following condition:

$$C = \left\{ p_i : L(p_i) > T_1 \quad \text{AND} \quad p_i = \underset{q \in b_i}{\text{argmax}} L(q) \right\}, \quad (3)$$

where  $p_i$  is an extracted feature point if its  $L(p_i)$  value is over a given threshold  $T_1$  and it is a local maximum in  $b_i$ . The  $T_1$  threshold is calculated adaptively by Otsu's method.<sup>19</sup> The  $C$  set of contour points is shown in Fig. 2(b). It is important to note, that MHEC also emphasizes parts that were dismissed by the original Harris implementation, like the frontal part of the left plane.

Now the  $C$  point set is defining contour points in the image belonging to different flying objects or background. The next step is to separate point subsets of various objects, while eliminating the points of the background.

The separation process of contour point subsets is based on the included points' connectivity in the Canny edge map.<sup>20</sup> If two contour points are connected by an edge in the edge map, then they are supposed to belong to the same object. The following graph representation formalizes this assumption: a  $G = (C, N)$  graph is described with the  $C$  vertex set and the  $N$  edge set, where  $C$  is defined by Eq. (3), and  $N$  is built according to the connectivity of the vertices (points) in the edge map.

Let  $E$  mark the dilated Canny edge map of the image [see Fig. 3(a)], where pixels representing the detected edges get a value of 1; others, representing the background, are 0. Two given vertices  $v_i, v_j \in C$  are connected by an edge in  $N$  if they satisfy the following conditions:

- (1)  $E(v_i) = 1$ ;  $v_i$  is an edge point in  $E$ .
- (2)  $E(v_j) = 1$ ;  $v_j$  is an edge point in  $E$ .
- (3) A finite path of pixels in  $E$  with value 1 exists between  $v_i$  and  $v_j$  (i.e., they are connected by an edge in  $E$ ).

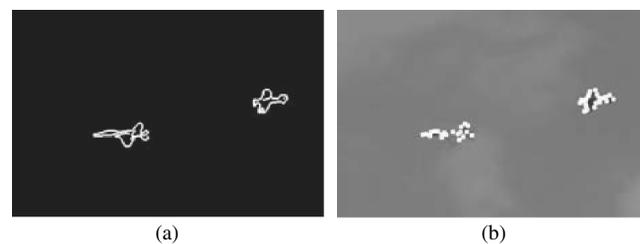
After performing this procedure for all vertices, the  $N$  edge set is defined.

Now the  $G$  graph will contain  $K$  disjoint subgraphs (denoting the  $k$ 'th with  $G^k$ ) with contour point sets  $C^k$  representing separate objects:

$$C^k = \{c_1^k, \dots, c_{N_k}^k\}, \quad (4)$$

where  $N_k$  is the number of contour points in  $G^k$ . Then the following conditions are satisfied by  $C^k$  point subsets:

$$C = \bigcup_{k=1}^K C^k; \quad C^i \cap C^j = \emptyset \quad \forall i, j. \quad (5)$$



**Fig. 3** Object separation: (a) Canny edge map; (b) separated object contour points.

Subgraphs containing only a few points are supposed to indicate noise or background, therefore we filter out a  $\mathbf{G}^k$  subgraph if the number of its points is smaller than an  $n$  threshold ( $N_k < n$ ). In our work, we applied  $n = 4$ . After this filtering, the remaining set of contour points representing  $K'$  flying targets is given as:

$$\mathbf{C}' = \bigcup_{k=1}^{K'} \mathbf{C}^k. \quad (6)$$

Figure 3(b) shows the two separated contour point sets representing the two objects.

We have to emphasize here, that while the MHEC point detection process produces more points than e.g., the Harris detector, this is a positive property of the approach, since we use and exploit the higher number of points to find good, connected boundary point locations in the graph analysis process. Thus, we arrive to a boundary point extraction, which is robust and at the same time aids in the removal on non-objects from the processed frames.

In situations where objects are so close that they visually occlude one another, the objects might not get separated and the result would be a blob containing both objects. In such situations the recognition phase (Sec. 3) will give a false classification of the blob. However, as it will be described later, the recognition step builds continuous statistics of the detected classes of objects in time, and if the objects will visually separate later, then their classes will be updated. Also, a tracker using the outputs of this paper could help in separation of such objects.

## 2.2 Object Extraction

In the previous step, we obtained corner and edge points which directly relate to the boundaries of flying objects. Having this information, we create a spatio-temporal, multimodal background model from the pixels which are not related to any of these objects. This background model is used for preliminary object extraction to roughly segment the objects' silhouettes from the background (see Sec. 2.2.1). Then, we create a separate appearance model for each object, and refine the silhouette in a Markov random field (MRF) framework using the background and the object appearance (i.e., foreground) models with an additional interaction constraint to express the similarity between neighboring pixels (see Sec. 2.2.2).

Formally, we denote by  $\mathbf{C}^k = \{c_1^k, \dots, c_{N_k}^k\}$  the set of corner and edge points (referred to as boundary points in the rest of the paper) of the  $k$ 'th object detected by the MHEC detector presented above, where  $N_k$  denotes the number of boundary points. Let  $\mathcal{C}_t = \{\mathbf{C}_t^1, \dots, \mathbf{C}_t^{K_t}\}$  denote the collection of all boundary point sets at time  $t$ , where  $K_t$  denotes the number of boundary point sets. Hereafter we denote by  $\mathbf{C}_t^k \in \mathcal{C}_t$  one element of the collection, where  $1 \leq k \leq K_t$ . Moreover, we also make the following formal definitions:

- $\mathbf{S}$  denotes the pixel lattice of an image;
- $\mathbf{X} = \{x_s | s \in \mathbf{S}\}$  is the set of pixel values of an image, i.e.,  $x_s$  is 3-tuple in a given color space;

- $L = \{b, f\}$  denotes the set of two class labels *background* (b) and *foreground* (i.e., object, f), respectively;
- $\Omega = \{\omega_s | s \in \mathbf{S}\}$  denotes the labeling of the pixels of an image, where  $\omega_s \in \mathbf{L}$  is the label of a particular pixel  $s$ ;
- $p_l(s) = P(x_s | \omega_s = l)$  is the conditional probability density function of label  $l \in \mathbf{L}$  at a given pixel  $s$ .

In our implementation we use the CIE  $L^*u^*v^*$  uniform color space, i.e.,  $x_s = [x_L(s), x_u(s), x_v(s)]$ .

### 2.2.1 Preliminary segmentation

Several pixel-level background estimation techniques exist,<sup>21,22</sup> however, these methods require a static camera to construct pixel-level statistical models, which makes them unfeasible for video sources with flying objects, where the camera is typically not static or sometimes even follows the moving airplane and the background can change from frame to frame, e.g., clouds or vapor trails might be visible, and illumination changes can also occur.

Because of the above problems, we create one global spatio-temporal background model at each timestep using the pixel values in a small moving time window. Here we use the pixels which do not relate to object silhouettes. In our method this global background  $p_b(s)$  is modeled with a finite mixture of  $M_b$  Gaussians (MoG), i.e.,

$$\begin{aligned} p_b(s) &= \sum_{k=1}^{M_b} \omega_{b,k} \cdot p_{b,k}(s) \\ &= \sum_{k=1}^{M_b} \omega_{b,k} \cdot \mathcal{N}(x_s | \mu_{b,k}, \Sigma_{b,k}), \end{aligned} \quad (7)$$

where  $p_{b,k}(s) = \mathcal{N}(x_s | \mu_{b,k}, \Sigma_{b,k})$  denotes the 3-D Gaussian density function, i.e.,

$$p_{b,k}(s) = \frac{\exp\left[-\frac{1}{2}(x_s - \mu_{b,k})^T \Sigma_{b,k}^{-1} (x_s - \mu_{b,k})\right]}{(2\pi)^{3/2} \cdot |\Sigma_{b,k}|^{1/2}}. \quad (8)$$

Moreover, to speed up the segmentation process we assume a diagonal covariance matrix, i.e.,  $\Sigma_{b,k} = \sigma_{b,k}^2 \cdot I$ , where  $I$  is the  $3 \times 3$  identity matrix.

At this point we utilize the  $\mathcal{C}_t$  collection of boundary point sets at time  $t$  as follows. Let  $b_t^k$  denote the bounding box of the  $k$ 'th point set  $\mathbf{C}_t^k \in \mathcal{C}_t$ , where the size of the box is slightly enlarged, and  $\mathbf{B}_t = \{b_t^1, \dots, b_t^{K_t}\}$  denotes the set of bounding boxes at time  $t$ . Moreover, let  $\mathbf{S}'_t \subset \mathbf{S}_t$  denote the pixels, which do not lie within any bounding boxes of  $\mathbf{B}_t$ , i.e.,

$$\mathbf{S}'_t = \{s \in \mathbf{S} : s \notin b_t^k, k = 1, \dots, K_t\}, \quad (9)$$

and let  $N_t$  denote the number of these pixels, i.e.,  $N_t = |\mathbf{S}'_t|$ . Furthermore, let  $r$  denote the radius of a small temporal window, and  $\mathbf{S}'_t(r)$  the union of pixels of all  $\mathbf{S}'_t$  sets in the  $r$  radius, i.e.,

$$\mathbf{S}'_t(r) = \bigcup_{\tau=-r}^{+r} \mathbf{S}'_{t+\tau}. \quad (10)$$

From Eqs. (9) to (10) it is obvious that

$$|\mathbf{S}'_t(r)| = N_t(r) = N_{t-r} + \dots + N_t + \dots + N_{t+r}. \quad (11)$$

Figure 4 demonstrates this process using one frame from a video with two objects (i.e.,  $K_t = 2$ ), where the inner rectangles (darker color) represent the  $b_t^1$  and  $b_t^2$  bounding boxes enlarged in both directions with 20%. The pixels outside these boxes are used to form the  $\mathbf{S}'_t$  set of pixels at time  $t$ . Having the  $\mathbf{S}'_t(r)$  set of training samples in the temporal window with radius  $r$  we calculate the maximum likelihood estimate of the global background MoG model using the expectation-maximization (EM) technique.<sup>23</sup>

The estimated background MoG model  $p_b(\cdot)$  can be directly used for foreground-background separation to determine the pixel labels of  $\Omega$ . However, in this case, the separation is based only on the background model with the risk that some parts of the object might get classified as background. Thus, we use this scheme for creating a preliminary classification only, and the results obtained are used for refining the segmentation by using a separate local appearance model for each object (see Sec. 2.2.2).

To obtain a preliminary object extraction, we use the estimated MoG background model of Eq. (7) to separate the foreground from the background. In our method, a pixel  $s$  is classified as background (i.e.,  $\omega_s = b$ ) in either of the following two cases:

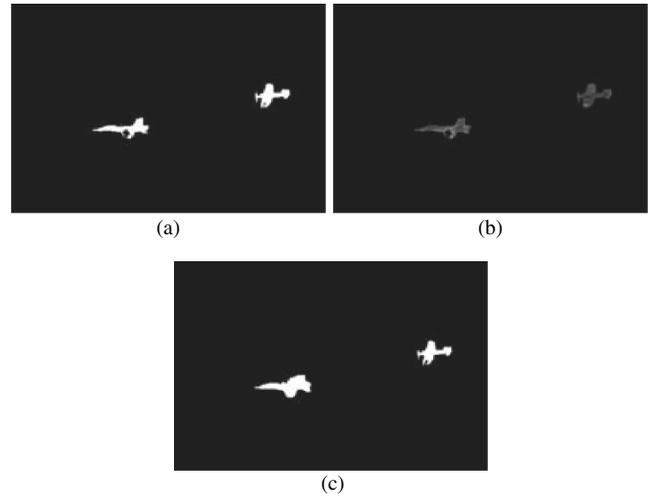
- (1) The pixel's position is "far" from the objects.
- (2) The pixel's value does not "match" the background model.

Otherwise the pixel is classified as foreground, i.e.,  $\omega_s = f$ . For the first case we simply use another rectangle around the boundary points, which has the size of the bounding box enlarged with 50% (illustrated by the outer rectangles with lighter color in Fig. 4), and we classify the pixels outside these rectangles as background. In the second case we considered pixel  $s$  "matching" the background model if for any Gaussian component of the mixture the following inequality holds:

$$\sqrt{(x_s - \mu_{b,k})^T \Sigma_{b,k}^{-1} (x_s - \mu_{b,k})} < T, \quad (12)$$



**Fig. 4** The inner (darker color) rectangles represent the bounding boxes of the detected edge and corner points. Pixels outside the inner rectangles are used for training the global background model. Pixels outside the outer rectangles (lighter color) are always classified as background in the preliminary segmentation step.



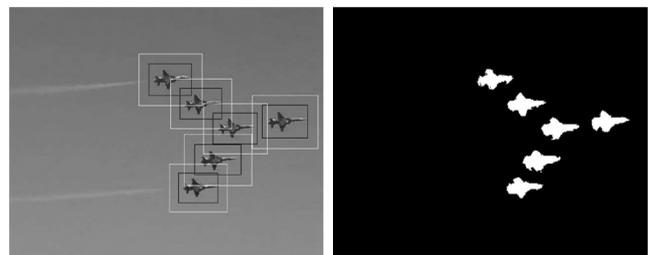
**Fig. 5** Foreground-background separation process of Fig. 4: (a) preliminary segmentation result; (b) pixel values used for estimating the foreground models; (c) final result.

where threshold  $T$  typically takes values on the [2.5;3.5] range, and we used a constant  $T = 2.5$  value in our experiments. The output of the preliminary segmentation of the Fig. 4 input is demonstrated in Fig. 5(a), where black color represents the pixels classified as background ( $\omega_s = b$ ), and white pixels denote the foreground pixels ( $\omega_s = f$ ). We can observe that a significant part of the right wing of the left airplane has been classified as background, since it matches the color of the clouds. This issue will be addressed in the next section by using local appearance models.

Situations might occur when two or more objects are so close to each other that their rectangle regions overlap, i.e., parts of one object become visible in the region of another object. Since the above described segmentation step is pixel-based and multimodal, in such situations the overlapping part of the other object will be extracted as a separate blob (will be classified as foreground, but as a separate region), and its remainder parts will become a different blob, extracted from a neighboring rectangular area. However, this causes no problems, since the output of the segmentation is a frame with white masks where it is not important whether a mask is a result of separate smaller blobs or was detected as a single larger blob. Figure 6 presents such a situation.

## 2.2.2 Markov random field segmentation

The preliminary segmentation process produces initial object silhouettes which might be broken or some parts of the objects may be misclassified as background. In the next step,



**Fig. 6** Situation where rectangles of objects overlap (left), and the produced mask.

we refine these silhouettes, and we follow a Bayesian approach for the classification of background and foreground pixels. For this step we need statistical information about the a priori and conditional probabilities of the two classes and the observable pixel values. In addition, we use an MRF to model the spatial interaction constraint of the neighboring pixels.

In our case the appearance of a foreground object is modeled by a MoG with the following conditional probability function:

$$p_f(s) = \sum_{k=1}^{M_f} \omega_{f,k} \cdot \mathcal{N}(x_s | \mu_{f,k}, \Sigma_{f,k}), \quad (13)$$

where we use a smaller number of components in the mixture than we used in the global background model, i.e.,  $M_f < M_b$ . To estimate the model parameters, we use the pixels within the  $b_f^k$  bounding box, which were classified as foreground in the preliminary segmentation step [highlighted in Fig. 5(b)]. The parameters of the MoG are obtained again by EM.<sup>23</sup>

According to the MRF model, the optimal labeling  $\hat{\Omega}$  is expressed as

$$\hat{\Omega} = \operatorname{argmin}_{\Omega} \sum_{s \in \mathcal{S}} -\log P(x_s | \omega_s) + \sum_{r,s \in \mathcal{S}} V(\omega_r, \omega_s), \quad (14)$$

where the spatial constraint is realized by the  $V(\omega_r, \omega_s)$  function (also known as smoothing term). Here we use  $V(\cdot, \cdot)$  to penalize those pixels whose class labels differ from those of their neighboring pixels. In our model  $V(\omega_r, \omega_s) = 0$  if  $\omega_r$  and  $\omega_s$  are not neighbors, otherwise

$$V(\omega_r, \omega_s) = \begin{cases} 0 & \text{if } \omega_r = \omega_s, \\ \beta & \text{if } \omega_r \neq \omega_s, \end{cases} \quad (15)$$

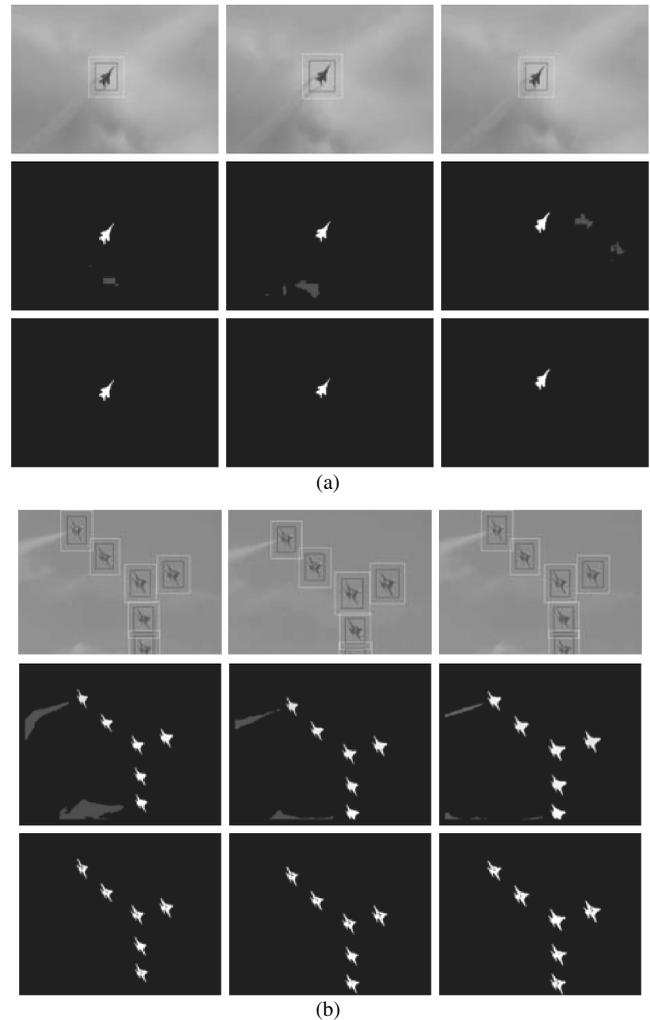
where  $\beta > 0$  is a penalizing constant.

Here we utilize the conditional probability functions of the two classes (background and foreground, respectively) at a given pixel  $s$ , defined in Eqs. (7) and (13). Finally, for solving the MRF problem we optimize the Eq. (14) energy function using a graph cuts based optimization algorithm,<sup>24–27</sup> which gives a good suboptimal solution in a few iteration of steps. Figure 5(c) shows the final result obtained by MRF segmentation. We can observe that the number of falsely classified pixels significantly decreased, e.g., the wings of the left airplane.

Figure 7 presents final outputs of the full presented object extraction approach above, compared to classical approaches,<sup>28</sup> which are less robust against background noise (e.g., clouds and vapor trails might be classified as foreground).

### 2.3 Feature Extraction

As later described, the recognition step of the proposed approach is not real time; thus, in our implementation, the queries against the shape dataset run in parallel with the main processing thread, updating the current object class after each step. A lower number of updates would generally mean lower recognition precision; thus, to alleviate this issue, we try to make the recognition itself more robust. We accomplish this by using a combination of two features—shape

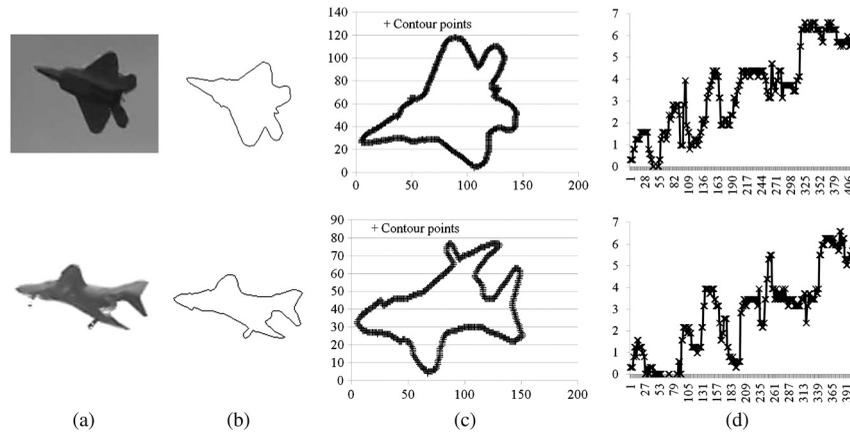


**Fig. 7** Samples for robustness against background noise (clouds, vapor trails, etc). For both (a) and (b), first row: input frame with background/foreground calculation regions (see Fig. 4); second row: classical foreground masks; third row: presented approach.

and texture—and we show that this solution helps in keeping (or even improving) the recognition accuracy even in the case of lower update frequencies. This in turn reduces the overall computational requirements of the whole process.

Shape features have been extracted and compared with a variety of methods in the literature, including hidden Markov models, scale invariant feature points, tangent/turning functions,<sup>29</sup> curvature maps, shock graphs, Fourier descriptors,<sup>30</sup> polar coordinates,<sup>31</sup> edge-based approaches like Chamfer distance methods (based on using small shape fragments),<sup>32</sup> and so on. They all have their benefits and drawbacks, regarding computational complexity, precision capabilities, implementation issues, robustness, and scalability. Overall, such methods convert some high-level description into a distance-based comparison, using some kind of chain code shape representation, incorporating rotation and scale invariance at the high level.

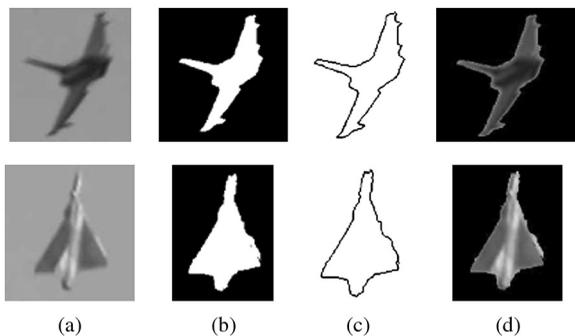
In our solution, we intended to use a shape description with a low computational complexity feature extraction step. Thus, we used a simple blob shape extraction, going over the obtained boundary points of the objects, and stored them as raw contour information. Also, we used a simple filtering



**Fig. 8** (a) Input frame object region; (b) object outline; (c) internal contour representation; and (d) turning function representation.

step to drop erroneous/noisy contour pixels (by a median-like filtering of the contour with a small neighborhood window), then used a scale and rotation invariant turning function representation for describing the contours, which are also be the basis for comparison (Sec. 3). For such a fairly simple method to work reliably, we need the outputs of the above presented robust foreground and object extraction step. Figure 8 shows some examples for this kind of object contour representations.

Texture features can also be extracted by a variety of methods, including LBP,<sup>33</sup> Gabor filters,<sup>34</sup> etc. To have a more complex texture representation with a tested method with known average performance, we decided to use the homogeneous texture descriptor (HTD)<sup>35</sup> from among the MPEG-7 descriptors. The HTD is a well-established and known standard texture descriptor, with known properties, extensive available tests and literature, and it is also lightweight. Reference 36 provides a detailed analysis and evaluation of the HTD against other features, as an effective way to describe object segmentation and image and video contents. The descriptor is computed by first filtering the image with a bank of orientation and scale sensitive Gabor filters, and computing the mean and standard deviation of the filtered outputs in the frequency domain. Figure 9 shows extracted shape examples from the previously obtained foregrounds and the respective regions that will be used for texture feature extraction.



**Fig. 9** (a) Section of the input frames with object; (b) extracted object blob; (c) extracted object contour; and (d) region of the original frame with texture.

### 3 Target Recognition

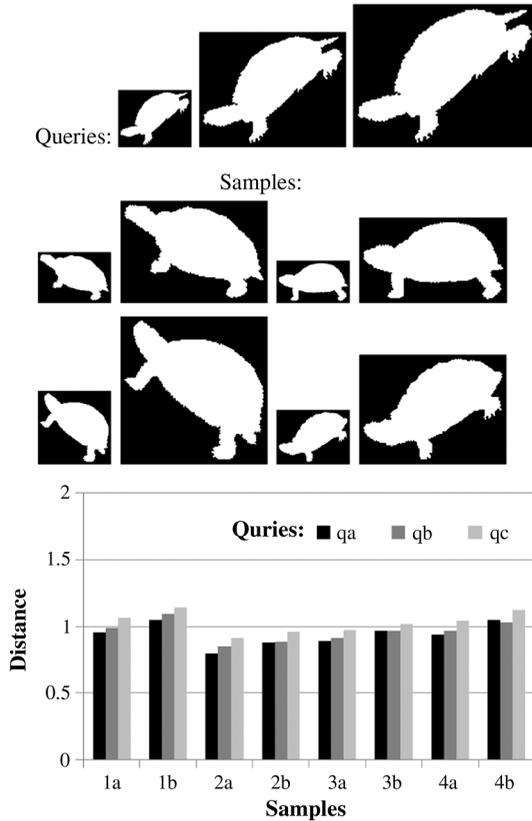
The goal of the recognition is to classify targets into one of the known classes during their observation, in a continuously updated process, regardless of the orientation, view, or size of the target. Generic object recognition methods are usually able to classify objects from the same view with which they have been trained. We do not want to have this limitation, so we build an index of target classes using various recordings containing targets moving freely, and the classification process is considered as a result produced for a query using this pre-built index. This approach enables us to quickly extend the recognizable classes, and to easily add more samples for a single class into the index, with the goal of recognizing the same target class from as many views as possible.

For comparing object shapes, we use the mentioned turning function representation,<sup>37</sup> mostly for speed and efficiency, where the object boundary is represented by a two-dimensional function describing the direction of the tangents of the curve along the objects' contours. To compare two such representations, the minimum distance of all shifted (by  $t$ ) and rotated (by  $\theta$ ) versions of the turning functions are produced (for rotation and scale invariant comparison), i.e.,

$$d_s^2(t_1, t_2) = \operatorname{argmin}_{t, \theta} \left[ \sum_k |t_1(k+t) - t_2(k) + \theta|^2 \right]. \quad (16)$$

The turning function comparison we use is based on the rotation- and scale-invariant method from Ref. 37, with the addition of a smoothing step to filter outlier points and sudden irregularities along a contour. Scale- and rotation-invariance of the function is also visualized in Fig. 10, where a query and its two scaled versions (top row, noted by  $qa$ ,  $qb$ , and  $qc$  in the bottom graph) are used as a basis for comparison against four other shapes which differ in orientation and in shape as well (differences at the head, leg and tail regions), and a scaled version of these four shapes (eight in total). Then, the graph shows the distance of each  $qa$ ,  $qb$ , and  $qc$  query against the rotated and scaled samples, showing that the bars in group are very similar to each other, meaning the differences of the scaled queries are very similar against all the rotated and scaled samples.

Texture features are compared by the HTD distance metric (from the MPEG-7 reference), comparing the local means and deviations obtained from the descriptor in the frequency domain. Distance is calculated as the minimum of

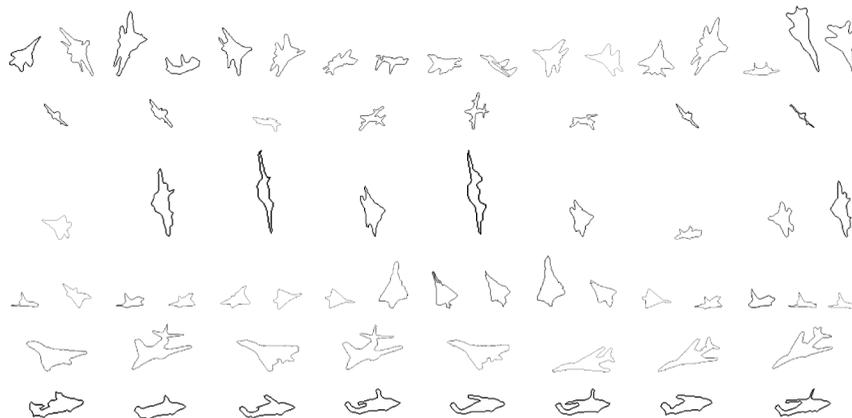


**Fig. 10** Visual samples for rotation- and scale-related properties of the used turning function-based calculations. The top row shows the three queries used (one query image and its two scaled versions, represented by  $q_a$ ,  $q_b$ , and  $q_c$  in the graph), while the samples are four objects with differences in shape, scale, and orientation. For each sample, the  $q_a$ ,  $q_b$ , and  $q_c$  bars are very similar, showing that the scaled and rotated samples are at the same distance from the scaled queries.

$L_1$  norms of shifted texture descriptor vectors ( $v_1$  and  $v_2$ ) for rotation invariance<sup>35</sup> as

$$d_t(v_1, v_2) = \min \left[ \sum_k |v_{1,m,f}(k) - v_2(k)| \right], \quad (17)$$

where  $f = 30^\circ$  is the angular division,  $m = 1, \dots, 5$ , and  $v_{1,m,f}$  is the shifted  $v_1$  vector. Experiments<sup>35</sup> have shown



**Fig. 11** Example for shapes from dataset classes (each line showing examples from a different class). Each class contains objects with different scale and orientation. Also a good example to show why such shape classes could not be easily learned by traditional classifiers.

an average accuracy of 77% for this descriptor (tested on the Brodatz texture dataset).

Since the class recognition step is not real time, queries against the indexed shape database are run at lower frequencies than the input video frame rate, and classification is done by building a continuous probability statistics of the results. As Fig. 1 shows, while processing of input frames is continuous (branch A), querying against the index runs in parallel (branch B), with reduced speed, updating the retrieval statistics after each cycle. In the following section dealing with evaluation of the recognition, we include details on how the combined recognition performs against changing query frequencies (i.e., how often queries are run against the index).

For recognition, an index is built from shape and texture features of a dataset of real-life videos, mostly from public recordings from air shows. The index structure is based on BK-trees.<sup>16</sup> Such index trees are representations of point distributions in discrete metric spaces. For classical string matching purposes, the tree is built so that each subtree contains sets of strings which are at the same distance from the subtree's root, i.e., for all  $e$  leaves below subroot  $r$ ,  $d(e, r) = \epsilon$  is constant. In our case, the used structure contains tree nodes that can have an arbitrary number of children ( $N$ ), where the leaves below each child contain elements for which the distance  $d$  falls in a difference interval:  $d(e, r) \in [\epsilon_i; \epsilon_{i+1})$ , where  $i \in [0, N] \cap \mathbb{N}$ . The distance intervals in the child nodes (denoted by  $\epsilon_i, \epsilon_{i+1}$  above) depend on the maximum error  $E_{\max}$  that the feature-dependent distance metric can have; more specifically  $\|\epsilon_{i+1} - \epsilon_i\| = E_{\max}/N$ , and thus, the difference intervals are linearly divided buckets.

The class recognition step is part of the main algorithm, but runs as a parallel process, using the outputs of the object extraction and feature extraction steps (branch B in Fig. 1). An important part of this process is the class candidate update, in which the recognition results are updated based on the statistics of the most probable results, and the best candidates are continuously refined based on the frequency of the result probabilities:

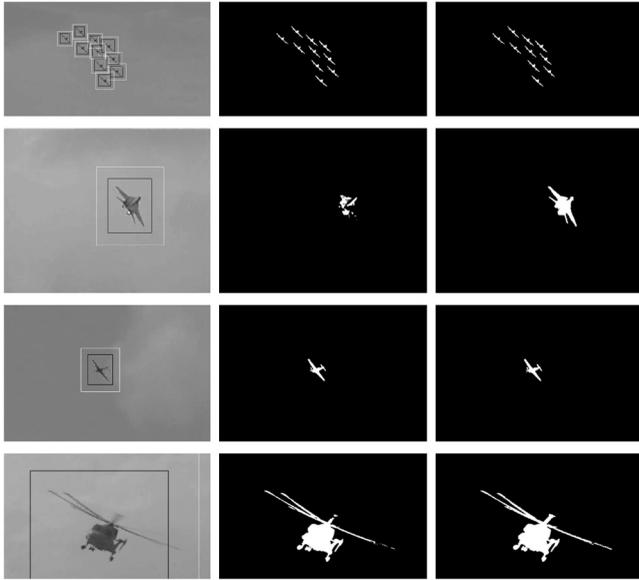
$$P_s(c_i) = P(c_i \in C_s) = \frac{\alpha \cdot nC_s}{\alpha \cdot nC_s + \beta \cdot nC_t},$$

$$P_t(c_i) = P(c_i \in C_t) = \frac{\beta \cdot nC_t}{\alpha \cdot nC_s + \beta \cdot nC_t}, \quad (18)$$

**Table 1** Lengths of query videos used in the test process.

Query video	1	2	3	4	5	6	7	8	9	10
Video length (frames)	322	134	67	55	100	132	130	143	75	220

where  $P_s$  and  $P_t$  are the probabilities of object  $c_i$  belonging to a specific shape ( $C_s$ ) and texture ( $C_t$ ) class, respectively,  $\alpha$  and  $\beta$  are weights that influence how the shape and texture information is included in the retrieval (typically shape has higher weight), and  $nC_s$ ,  $nC_t$  is how many times the object



**Fig. 12** Example processing outputs for detection of objects. For each example, left to right: input frame with regions for background modeling; raw foreground masks; final filtered object masks.

has been classified in to the same specific class during the observation of the object during the full input video (or camera stream), i.e.,:

$$nC_s(c_i) = \sum_{k=1}^{i-1} P[c_k \in C_s(c_i)], \quad (19)$$

where  $c_i$  is the currently evaluated target, and the sum gathers all the instances where this target belonged to the same  $C_s$  class throughout the observation period.

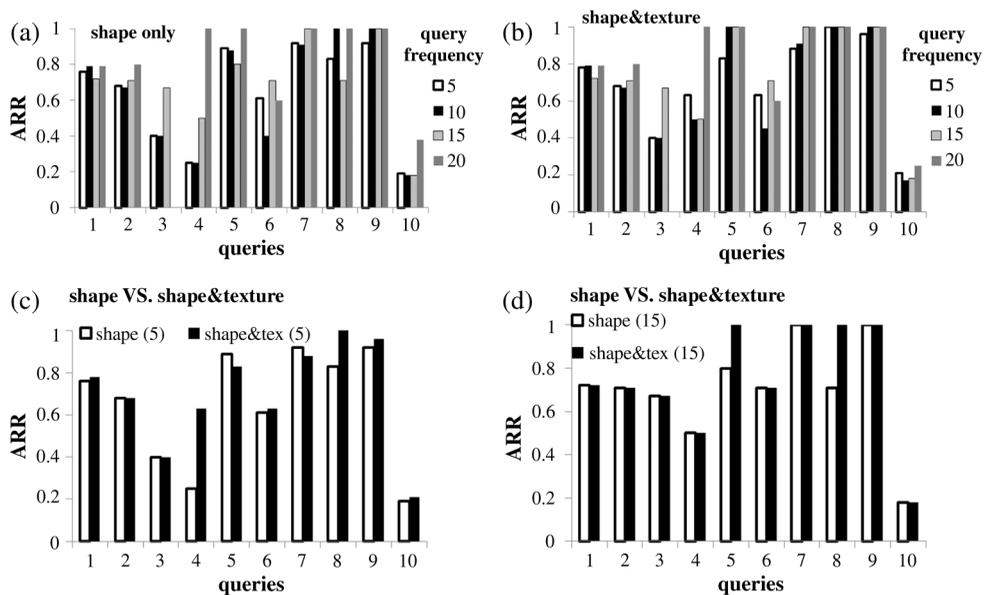
The probability of object  $c$  belonging to a specific class  $C$  identified by the respective  $C_s, C_t$  shape and texture classes will be the class into which it has been observed the most frequently during the observation period:

$$P(c \in C) = \max\{P_s, P_t\}. \quad (20)$$

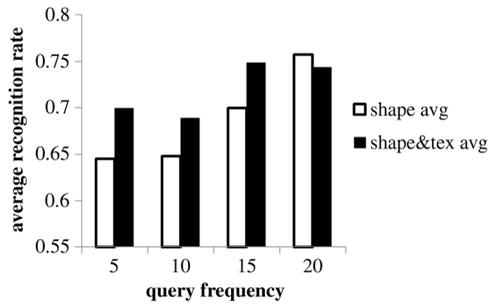
### 3.1 Data and Evaluation

For evaluation, we used a dataset gathered from public real life air show videos. The dataset contains 26 classes of planes, and the shapes and textures were extracted automatically with the above-presented methods. For each dataset video, objects have been extracted from all frames, which results in classes including a very high variation in scale and orientation of the specific targets (e.g., Fig. 11 shows examples from classes to illustrate typical intra-class diversity).

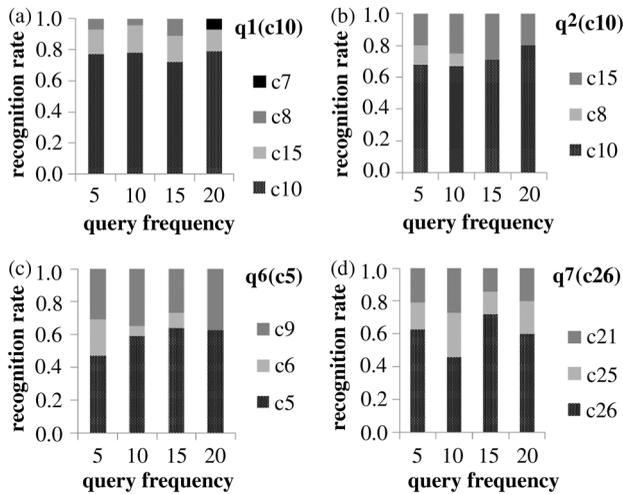
Each class has been manually labeled, and indexes have been built from the extracted shape and texture data. The indexing process for the approximately 8500 objects takes



**Fig. 13** Evaluation graphs for recognition rates: (a) and (b) show shape-only and combined shape + texture accumulated recognition rates (ARR) over 10 query videos and four different query frequency settings; (c) and (d) show shape-only versus combined shape + texture recognition rates in two selected frequency settings (every five and 15 frames).



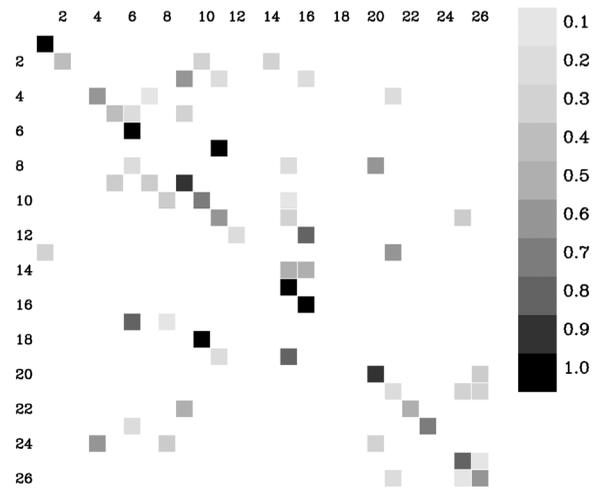
**Fig. 14** Averaged recognition rates (over all query videos) for shape only and combined queries, for different query frequencies.



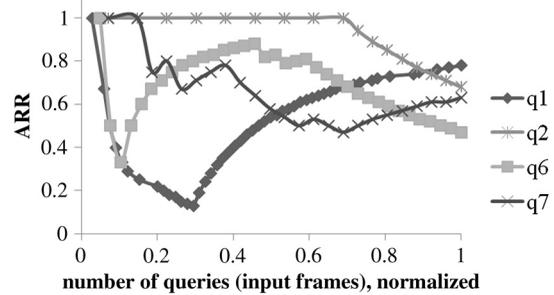
**Fig. 15** Evaluation graphs showing accumulated recognition rates for different queries (1, 2, 6, and 7) for different query frequencies. The columns show for a given query frequency which are the first three best matches during the recognition process. For example, a 0.8 recognition rate means that during the entire recognition process, the given specific class was given as the recognition result in 80% of the queries.

**Table 2** Samples from query video frames and samples from the best matches at the end of the query video (according to the accumulated recognition rates).

Query	Input frames	Best matches
q1		✓ 79% 17% 11%
q2		✓ 71% 29%
q6		✓ 64% 27% 9%
q7		✓ 71% 14% 14%



**Fig. 16** Confusion matrix of the used classes, showing color coded values of recognition rates, using queries outside of the indexed dataset. Numbers along the axes represent class numbers.



**Fig. 17** Accumulated recognition rates versus the length of the recognition process. The horizontal axis represents the length of the input videos (i.e., the numbers of different queries) on which recognition is being performed, all normalized to one for easier visualization.

approximately 2 min. on a 2.8-GHz Intel Core i7™ CPU, but this is done off-line, and the produced indexes are used during the actual processing (see Fig. 1).

To evaluate the retrieval/classification, we used videos not part of the dataset for testing, that contained objects from classes of the dataset. Table 1 shows the length (in frames) of the used testing videos.

Regarding the the object extraction discussed in Sec. 2.2, we used the following settings. The number of components of the MoG background and foreground were set to  $M_b = 6$  and  $M_f = 3$ , the temporal window had a radius  $r = 5$ . In the MRF model we used  $\beta = 1/2$ , and to optimize the energy of the model we used Szeliski et al.'s MRF minimization.<sup>27</sup> Examples of generated object masks are presented in Fig. 12, containing samples for multiple objects and challenging contents (clouds, illumination changes, etc.). Sources for the object extraction and some video examples are available online (<http://web.eee.sztaki.hu/~ucu/sw/>).

For evaluating the recognition rates of the proposed scheme, queries with different frequencies (i.e., every 5th, 10th, 15th, and 20th frames) were performed on each of the query videos, with a total of 1,118 retrievals, producing the results in Figs. 13 and 14.

The reason for testing the target classification process with varying query frequencies is because the querying process is not real time, recognition runs as a parallel process

**Table 3** Recognition rates from other works dealing with classification of synthetic or manually segmented shape classes.

Ref. 38—avg. precisions for:		Ref. 30—avg. precisions for:		Ref. 31—avg. precisions for diff. features:	
Different noise levels:	0.92	99 shapes	0.9	MI	0.69
30 degree tilt for 9 slant levels:	0.8	216 shapes	0.97	FD	0.57
60 degree tilt for 9 slant levels:	0.75	1045 shapes	0.84	UNL	0.71
90 degree tilt for 9 slant levels:	0.69	24 shape classes	0.9	UNL-F	0.98

(see Fig. 1), and to reduce the overall computational requirements, it would be desirable to run a target query process less frequently. The goal of the tests is to show that by using shape- and texture-fused recognition, reducing the query frequency does not hinder the classification process, because errors eventually caused by the less frequent update of the statistics are balanced by the more robust classification using the combined features.

Figure 13 contains results showing recognition rate data for the 10 query videos, for five different query frequencies. Figure 13(a) and 13(b) shows recognition rates for shape-only evaluation (only shape-based queries were performed) and combined shape + texture evaluation. These graphs show that the combined recognition is less dependent on the query frequency as when only using shape; moreover, the combined feature rate is higher for most of the queries. Figure 13(c) and 13(d) shows comparisons of rates between shape-only and combined shape + texture queries for two separate query frequencies (every five and 15 frames, respectively). The graphs show that the combined recognition rates are typically as good or better than the shape-only rates. This is an important feature, since high similarity in texture, which is very typical for planes, should cause a decrease in rates, but the weighted query scheme and the continuously updated class probabilities provide a balanced solution for these problems.

Figure 14 shows averaged recognition rates, where for each query frequency, the rate averages were computed over all 10 query videos, for both shape-only and combined shape + texture queries, which represents hundreds of averaged queries for each video. This graph shows that, on average, less frequent querying introduces less errors/noise in the recognition process, but also, that a combined shape + texture recognition scheme is generally better suited for such classification tasks than relying on shape information alone.

Figure 15 shows accumulated recognition rates (average precision values) for different queries and query frequencies, also showing the first best matches besides the recognized category. These graphs present in a visual form lines of so-called confusion matrices, showing how the recognition rates are distributed among the known classes (i.e., which classes and in what percentage are given as results for a query). As described earlier, the recognized class is produced as the best performing category, i.e., the result with the currently highest precision rate (these are the bottom parts of all columns in the figure). The figure's columns also show

the other top guesses, which follow the top result. Connected to this figure, Table 2 visually shows the results for different queries, where frames are sampled from the query and test videos. The “input frames” column shows samples from the query video, while the “best matches” column show the recognized result class and samples from other classes which are close to the result. Here, we can see even the second- and third-ranked classes are visually close to the correct response (which is the cause for lower precision values).

Connected to the previous figure, Fig. 16 shows the confusion matrix for all classes. Here, 10 to 15 frame-long queries were used for each class, which were not part of the indexed dataset, and recognition rates were recorded for each query, which are color coded in the figure. We must note here that as described above, recognition rates depend on a number of factors, e.g., length of the query, or number of examples that were included in the dataset index. As a general rule, recognition rates can be improved by observing the target for longer time periods, and if needed, by extending the dataset with the new sample and regenerating the index.

Figure 17 (in connection with Fig. 15) shows the evolution of recognition rates in relation to the number of queries, i.e., how many times a target is attempted to be recognized. This in practice translates to the length of the observation period of a target, during which recognition (query-retrieval step) is performed with a specific frequency (as described earlier). As previous data also shows, the average recognition rates tend to converge around 75% in time (in the figure, the length of the observation period is normalized, for better visualization).

Other works in the field of shape recognition typically deal with already available datasets containing clear contours and no noise, or synthetically added noise/distortions. However, as a comparison with other approaches for shape recognition and retrieval, we have included some other results in Table 3. Our average recognition rate of 75% (Fig. 14) is amongst the averages of other approaches. The best advantage of the presented approach is that it produces the presented results on real data, including all steps from background modeling, shape extraction and recognition, and it is easily expandable, e.g., by adding a tracking algorithm on top of the produced results.

## 4 Conclusions

We presented a flying target detection method based on corner and edge detection combined with robust background

modeling, serving as the basis for a target recognition scheme, which uses a fusion of shape and texture features for indexing and classifying the extracted objects. The method's goals are to be lightweight and robust, providing solutions suitable for application in real-time visual systems for defensive surveillance scenarios, as a basis for passive sensors. The detected object contours, main corner and boundary points, and object features can be used for target recognition and tracking. Sample sources for object extraction and video examples are available online (<http://web.eee.sztaki.hu/~ucu/sw/>). Future work includes adaptation of the methods for ground object detection and tracking for more generic object classes, and adaptation of the presented algorithms for smart cameras in order to provide an integrated solution.

### Acknowledgments

The work has been partially supported by Grant No. 83438 of the Hungarian Scientific Research Fund (OTKA).

### References

- J. W. Lu, Y. J. He, and H. Y. Li, "Detecting small target of ship at sea by infrared image," in *Proc. IEEE Intl. Conf. Autom. Sci. Eng.*, pp. 165–169, IEEE, Shanghai (2006).
- S. Zheng, J. Liu, and J. W. Tian, "An SVM-based small target segmentation and clustering approach," in *Proc. Intl. Conf. Mach. Learn. Cybern.*, Vol. 6, pp. 3318–3323 (2004).
- H. Deng and J. Liu, "Infrared small target detection based on the self-information map," *Inf. Phys. Technol.* **54**(2), 100–107 (2011).
- L. Itti, C. Gold, and C. Koch, "Visual attention and target detection in cluttered natural scenes," *Opt. Eng.* **40**(9), 1784–1793 (2001).
- X. Lia et al., "Detecting missile-like flying target from a distance in sequence images," *Proc. SPIE* **6968** (2008).
- S. Zhang, W. Liu, and X. Xue, "Research on tracking approach to low-flying weak small target near the sea," *Proc. SPIE* **6027**, 962–968 (2006).
- C. Bibby and I. Reid, "Visual tracking at sea," in *Proc. IEEE Intl. Conf. Robotics Autom.* (2005).
- T. L. Wenga et al., "Weather-adaptive flying target detection and tracking from infrared video sequences," *Expert Sys. Appl.* **37**(2), 1666–1675 (2010).
- A. L. Chan, "A robust target tracking algorithm for FLIR imagery," *Proc. SPIE* **7696**, 769603 (2010).
- Y. Wang et al., "Multi-class target recognition based on adaptive feature selection," *Proc. SPIE* **7696**, 769609 (2010).
- E. Blasch and B. Kahler, "Multiresolution EO/IR target tracking and identification," in *Proc. Intl. Conf. Inform. Fus.* Vol. 1, pp. 275–282 (2005).
- H. Noor et al., "Model generation for video-based object recognition," in *Proc. ACM Intl. Conf. Multimedia*, pp. 715–719, ACM, New York, NY (2006).
- D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vis.* **60**(2), 91–110 (2004).
- L. Kovács and T. Szirányi, "Recognition of hidden pattern with background," *Proc. SPIE* **6699**, 669906 (2007).
- A. Kovács et al., "Shape and texture fused recognition of flying targets," *Proc. SPIE* **8050**, 80501E (2011).
- W. A. Burkhard and R. M. Keller, "Some approaches to best-match file searching," *Commun. ACM* **16**(4), 230–236 (1973).
- C. Harris and M. Stephens, "A combined corner and edge detector," *Proc. 4th Alvey Vis. Conf.*, pp. 147–151 (1988).
- A. Kovács and T. Szirányi, "Harris function based active contour external force for image segmentation," *Patt. Recogn. Lett.* **33**(9), 1180–1187 (2012).
- N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Sys. Man Cybern.* **9**(1), 62–66 (1979).
- J. Canny, "A computational approach to edge detection," *IEEE Trans. Patt. Anal. Mach. Intell.* **PAMI-8**(6), 679–698 (1986).
- K. Toyama et al., "Wallflower: principles and practice of background maintenance," in *Proc. 7th IEEE Int. Conf. Comp. Vis.*, Kerkyra, Greece, **1**, 255–261 (1999).
- C. Stauffer and W. E. L. Grimson, "Learning patterns of activity using real-time tracking," *IEEE Trans. Patt. Anal. Mach. Intell.* **22**(8), 747–757 (2000).
- A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *J. Royal Stat. Soc. B* **39**(1), 1–38 (1977).
- Y. Boykov, O. Veksler, and R. Zabih, "Fast approximate energy minimization via graph cuts," *IEEE Trans. Patt. Anal. Mach. Intell.* **23**(11), 1222–1239 (2001).
- V. Kolmogorov and R. Zabih, "What energy functions can be minimized via graph cuts?," *IEEE Trans. Patt. Anal. Mach. Intell.* **26**(2), 147–159 (2004).
- Y. Boykov and V. Kolmogorov, "An experimental comparison of min-cut/max-flow algorithms for energy minimization in vision," *IEEE Trans. Patt. Anal. Mach. Intell.* **26**(9), 1124–1137 (2004).
- R. Szeliski et al., "A comparative study of energy minimization methods for Markov random fields," in *Proc. 9th Eur. Conf. Comp. Vis.*, Graz, Austria, pp. 16–29 (2006).
- L. Kovács and Á. Utasi, "Shape and motion fused multiple flying target recognition and tracking," *Proc. SPIE* **7696**, 769605 (2010).
- L. J. Latecki and R. Lakamper, "Application of planar shape comparison to object retrieval in image databases," *Intl. J. Info. Sci.* **35**(1), 15–29 (2002).
- W. T. Wong, F. Y. Shih, and J. Liu, "Shape-based image retrieval using support vector machines, Fourier descriptors and self-organizing maps," *Intl. J. Info. Sci.* **177**(8), 1878–1891 (2007).
- D. Frejlichowski, "An algorithm for binary contour objects representation and recognition," in *Proc. ICIAR, Lecture Notes Comp. Sci.*, pp. 537–546 (2008).
- J. Shotton, A. Blake, and R. Cipolla, "Multi-scale categorical object recognition using contour fragments," *IEEE Trans. Patt. Anal. Mach. Intell.* **30**(7), 1270–1281 (2008).
- T. Ojala, M. Pietikäinen, and T. Maenpää, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Trans. Patt. Anal. Mach. Intell.* **24**(7), 971–987 (2002).
- A. K. Jain and F. Farrokhnia, "Unsupervised texture segmentation using Gabor filters," *Patt. Recogn.* **24**(12), 1167–1186 (1991).
- B. S. Manjunath et al., "Color and texture descriptors," *IEEE Trans. Circuits Sys. Vid. Technol.* **11**(6), 703–715 (2001).
- Y. M. Ro et al., "MPEG-7 homogeneous texture descriptor," *ETRI J.* **23**(2), 41–51 (2001).
- E. W. Arkin et al., "An efficiently computable metric for comparing polygonal shapes," *IEEE Trans. Patt. Anal. Mach. Intell.* **13**(3), 209–216 (1991).
- M. Bicego and V. Murino, "Investigating Hidden Markov Models' capabilities in 2D shape classification," *IEEE Trans. Patt. Recogn. Mach. Intell.* **26**(2), 281–286 (2004).



**Levente Kovács** is a senior research fellow at the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Budapest, Hungary. He has received his MSc in IT (2002) and PhD in image processing and graphics (2007) from the University of Pannonia, Hungary. His main research areas are image/video feature extraction and indexing, annotation, event detection, nonphotorealistic rendering, and video restoration.



**Andrea Kovács** is a PhD student and researcher at the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Budapest, Hungary. She received her MSc degree in electrical engineering from the Technical University of Budapest in 2008. Her research interests include image and video processing for remote sensing, shape analysis, boundary extraction, and active contours.



**Ákos Utasi** is a research fellow at the Computer and Automation Research Institute of the Hungarian Academy of Sciences, Budapest, Hungary. He received his MSc in CS (2005) and PhD in visual surveillance (2012) from the University of Pannonia, Hungary. His research interests include visual surveillance, motion detection, event detection, and action recognition.



**Tamás Szirányi** received his PhD and DSci degrees in 1991 and 2001, respectively, from the Hungarian Academy of Sciences. He was appointed as a full professor at Pannon University, Veszprém, Hungary, in 2001 and, in 2004, at the Péter Pázmány Catholic University, Budapest. He leads the Distributed Events Analysis Research Laboratory at the Computer and Automation Research Institute, Budapest. His research areas include machine perception, stochastic optimization, remote sensing, surveillance systems for multiview camera

systems, intelligent networked sensor systems, graph based clustering, and film restoration. He was the founder and past president (1997 to 2002) of the Hungarian Image Processing and Pattern Recognition Society, associate editor (AE) of IEEE Trans. Image Processing (2003 to 2009), AE of Digital Signal Processing since 2012. He was honored with the Master Professor (2001), and ProScientia (2011) awards. He is the senior member of IEEE and a fellow of the IAPR and the Hungarian Academy of Engineering. He has more than 200 publications, including 40 in major scientific journals.