



Real-time human segmentation in infrared videos

Antonio Fernández-Caballero^{a,b,*}, José Carlos Castillo^a, Juan Serrano-Cuerda^a,
Saturnino Maldonado-Bascón^c

^a Instituto de Investigación en Informática de Albacete (I3A), Universidad de Castilla-La Mancha, 02071-Albacete, Spain

^b Departamento de Sistemas Informáticos, Escuela de Ingenieros Industriales de Albacete, Universidad de Castilla-La Mancha, 02071-Albacete, Spain

^c Department of Signal Theory and Communications, Escuela Politécnica Superior, Universidad de Alcalá, 28871-Alcalá de Henares, Madrid, Spain

ARTICLE INFO

Keywords:

Human segmentation
Infrared image
Real-time performance

ABSTRACT

In this paper, a new approach to real-time people segmentation through processing images captured by an infrared camera is introduced. The approach starts detecting human candidate blobs processed through traditional image thresholding techniques. Afterwards, the blobs are refined with the objective of validating the content of each blob. The question to be solved is if each blob contains one single human candidate or more than one. If the blob contains more than one possible human, the blob is divided to fit each new candidate in height and width.

© 2010 Elsevier Ltd. All rights reserved.

1. Introduction

The use of infrared cameras in the surveillance field (López, Fernández-Caballero, Fernández, Mira, & Delgado, 2006; Pavón, Gómez-Sanz, Fernández-Caballero, & Valencia-Jiménez, 2007) are being intensively studied in the last decades. Many algorithms focusing specifically on the thermal domain have been explored. The unifying assumption in most of these methods is the belief that the objects of interest are warmer than their surroundings (Yilmaz, Shafique, & Shah, 2003). Thermal infrared video cameras detect relative differences in the amount of thermal energy emitted/reflected from objects in the scene. As long as the thermal properties of a foreground object are slightly different (higher or lower) from the background radiation, the corresponding region in a thermal image appears at a contrast from the environment.

In Iwasawa, Ebihara, Ohya, and Morishima (1997) and Bhanu and Han (2002), a thresholded thermal image forms the first stage of processing after which methods for pose estimation and gait analysis are explored. In Nanda and Davis (2002), a simple intensity threshold is employed and followed by a probabilistic template. A similar approach using Support Vector Machines is reported in Xu, Liu, and Fujimura (2005). Recently, a new background-subtraction technique to robustly extract foreground objects in thermal video under different environmental conditions has been presented (Davis & Sharma, 2007). A recent paper (Jung, Eledath, Johansson, & Mathevon, 2007) presents a real-time

egomotion estimation scheme that is specifically designed for measuring vehicle motion from a monocular infrared image sequence at night time. In the robotics field, a new type of infrared sensor is described (Benet, Blanes, Simó, & Pérez, 2002). It is suitable for distance estimation and map building. Another application using low-cost infrared sensors for computing the distance to an unknown planar surface and, at the same time, estimating the material of the surface has been described (García & Solanas, 2004).

In this paper, we introduce our approach to real-time robust people segmentation through processing video images captured by an infrared camera, complementing this way the previous works on knowledge-based object detection (Fernández-Caballero, Gómez, & López-López, 2008; Fernández-Caballero, López, & Saiz-Valverde, 2008; Fernández-Caballero et al., 2007; López, Fernández-Caballero, Mira, Delgado, & Fernández, 2006; Mira, Delgado, Fernández-Caballero, & Fernández, 2004).

2. Robust people segmentation algorithm

The proposed human detection algorithm is explained in detail in the following sections related to the different phases, namely, people candidate blobs detection, people candidate blobs refinement and people confirmation.

2.1. People candidate blobs detection

The algorithm starts with the analysis of input image, $I(x,y)$, captured at time t by an infrared camera, as shown in Fig. 1(a). Firstly, a change in scale, as shown in Eq. (1) is performed. The idea is to normalize all images to always work with a similar

* Corresponding author at: Instituto de Investigación en Informática de Albacete (I3A), Universidad de Castilla-La Mancha, 02071-Albacete, Spain. Tel.: +34 967 599200; fax: +34 967 599224.

E-mail address: caballer@dsi.uclm.es (A. Fernández-Caballero).

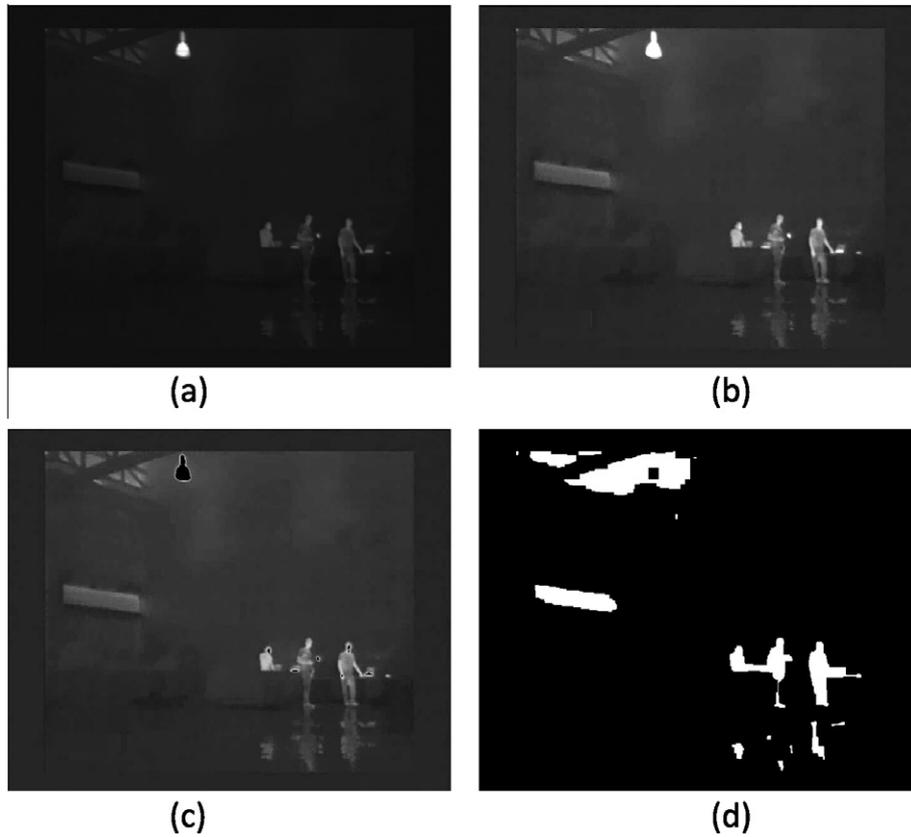


Fig. 1. (a) Input infrared image, (b) Scaled frame, (c) Incandescence elimination, (d) Thresholded frame.

scale of values, transforming $I(x,y)$ to $I^n(x,y)$ (see Fig. 1(b)). The normalization assumes a factor γ , calculated as the mean gray level value of the las n input image, and uses the mean gray level value of the current image, \bar{I} .

$$I^n(x,y) = \frac{I(x,y) \times \gamma}{\bar{I}} \quad (1)$$

where $I^n(x,y)$ is the normalized image. Notice that $I^n(x,y) = I(x,y)$ when $\bar{I} = \gamma$.

The next step is the elimination of incandescent points-corresponding to light bulbs, power sources, and so on-, which can confuse the algorithm by showing zones with too high temperatures. As the image has been scaled, the threshold θ_i calculated to eliminate these points is related to the normalization factor γ . Indeed,

$$\theta_i = 3 \times \frac{5}{4} \gamma \quad (2)$$

$\delta = \frac{5}{4} \gamma$ introduces a tolerance value of a 25% above the mean image value. And, $3 \times \delta$ provides a value high enough to be considered an incandescent image pixel. Thus, pixels with a higher gray value are discarded and filled up with the mean gray level of the image.

$$I^n(x,y) = \begin{cases} I^n(x,y), & \text{if } I^n(x,y) \leq \theta_i \\ \bar{I}^n, & \text{otherwise} \end{cases} \quad (3)$$

The algorithm uses a threshold to perform a binarization for the aim of isolating the human candidates spots. The threshold θ_c , obtains the image areas containing moderate heat blobs, and, therefore, belonging to human candidates. Thus, warmer zones of the image are isolated where humans could be present. The threshold is calculated as:

$$\theta_c = \frac{5}{4} (\gamma + \sigma^n) \quad (4)$$

where σ^n is the standard deviation of image $I^n(x,y)$. Notice, again, that a tolerance value of a 25% above the sum of the mean image gray level value and the image gray level value standard deviation is offered.

Now, image $I^n(x,y)$ is binarized using the obtained threshold θ_c . Pixels above the threshold are set as maximum value $max = 255$ and pixels below are set as minimum value $min = 0$.

$$I_b^n(x,y) = \begin{cases} min, & \text{if } I^n(x,y) \leq \theta_c \\ max, & \text{otherwise} \end{cases} \quad (5)$$

Next, the algorithm performs morphological opening (Eq. (6)) and closing (Eq. (7)) operations to eliminate isolated pixels and to unite areas split during the binarization. These operations require structuring elements that in both cases are 3×3 square matrices centered at position (1,1). These operations greatly improve the binarized shapes as shown in Fig. 1.

$$I_o^n(x,y) = I_b^n(x,y) \circ \begin{vmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{vmatrix} \quad (6)$$

$$I_c^n(x,y) = I_o^n(x,y) \bullet \begin{vmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{vmatrix} \quad (7)$$

Afterwards, the blobs contained in the image are obtained. A minimum area, A_{min} , – function of the image size- is established for a blob to be considered to contain humans.

$$A_{min} = 0.0025 \times (r \times c) \quad (8)$$

where r and c are the number of rows and columns, respectively of input image $I(x,y)$. As a result, the list of blobs, L_B , containing people candidates in form of blobs $b_\lambda[(x_{start},y_{start}),(x_{end},y_{end})]$, is generated. λ stands for the number (index) of people candidate blob in image $I^m(x,y)$, whereas (x_{start},y_{start}) and (x_{end},y_{end}) are the upper left and lower right coordinates, respectively, of the minimum rectangle containing the blob. As an example, consider the resulting list of blobs related to Fig. 1 and offered in Table 1.

2.2. People candidate blobs refinement

In this part, the algorithm works with the list of blobs L_B , present in image $I^m(x,y)$, obtained at the very beginning of the previous section. At this point, there is a need to validate the content of each blob to find out if it contains one single human candidate or more than one. Therefore, the algorithm processes each detected blob separately.

Let us define a region of interest (ROI) as the minimum rectangle containing one blob of list L_B (obtained from $I^m(x,y)$). A ROI may be defined as $R_\lambda = R_\lambda(i,j)$, when associated to blob $b_\lambda[(x_{start},y_{start}),(x_{end},y_{end})]$. Notice that $i \in [1 \dots max_i = x_{end} - x_{start} + 1]$ and $j \in [1 \dots max_j = y_{end} - y_{start} + 1]$.

2.2.1. People vertical delimiting

The first step consists in scanning R_λ by columns, adding the gray level value corresponding to each pixel in that column, as shown in Eq. (9).

$$H_\lambda[i] = \sum_{j=1}^{max_j} R_\lambda(i,j), \quad \forall i \in [1 \dots max_i] \tag{9}$$

This way, a histogram $H_\lambda[i]$ showing which zones of the ROI own greater heat concentrations is obtained. A double purpose is pursued when computing the histogram. In first place, we want to in-

crease the certainty of the presence and situation of human heads. Secondly, as a ROI may contain several persons that are close enough to each other, the histogram helps separating human groups (if any) into single humans. This method, when looking for maximums and minimums within the histogram allows differentiating among the humans present in the particular ROI.

Now the histogram, $H_\lambda[i]$, is scanned to separate grouped humans, if any. For this purpose, local maxima and local minima are searched in the histogram to establish the different heat sources (see Fig. 2(a)). To assess whether a histogram column contains a local maximum or minimum, a couple of thresholds are fixed, $\theta_{v_{max}}$ and $\theta_{v_{min}}$. Experimentally, we went to the conclusion that the best thresholds should be calculated as:

$$\theta_{v_{max}} = 2 \times \bar{R}_\lambda + \sigma_{R_\lambda} \tag{10}$$

$$\theta_{v_{min}} = 0.9 \times \bar{R}_\lambda \times max_j \tag{11}$$

Each different region that surpasses $\theta_{v_{max}}$ is supposed to contain one single human head, as heads are normally warmer than the rest of the people body covered by clothes. That is why $\theta_{v_{max}}$ has been set to the double of the sum of the average gray level plus the standard deviation of the ROI. On the other hand, $\theta_{v_{min}}$ indicates those regions of the ROI where the sum of the heat sources are really low. These regions are supposed to belong to gaps between two humans. We are looking for regions where the column summed gray level is below a 90% of the mean ROI gray level value. Fig. 2(b) shows the histogram for input ROI of Fig. 2(a). You may observe the values for $\theta_{v_{max}}$ and $\theta_{v_{min}}$, corresponding to three peaks (three heads) and two valleys (two separation zones). Fig. 2(c) shows the three humans as separated by the algorithm into sub-ROIs, $sR_{\lambda,\alpha}$.

2.2.2. People horizontal delimiting

All humans contained in a sub-ROI, $sR_{\lambda,\alpha}$, obtained in the previous section possess the same height, namely the height of the original ROI. Now, we want to fit the height of each sub-ROI to the real

Table 1
People candidates blobs list.

λ	x_{start}	y_{start}	x_{end}	y_{end}	Area
1	297	270	482	458	35154
2	608	344	645	376	1254

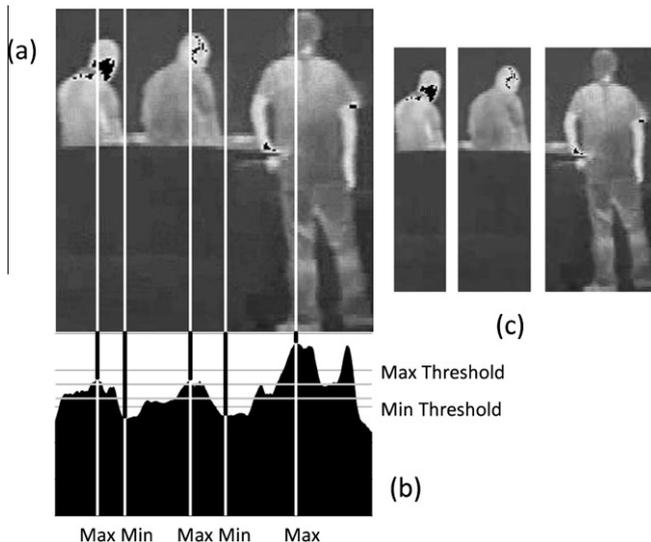


Fig. 2. (a) Input ROI, (b) Histogram, (c) Columns adjustment to obtain three human candidates.

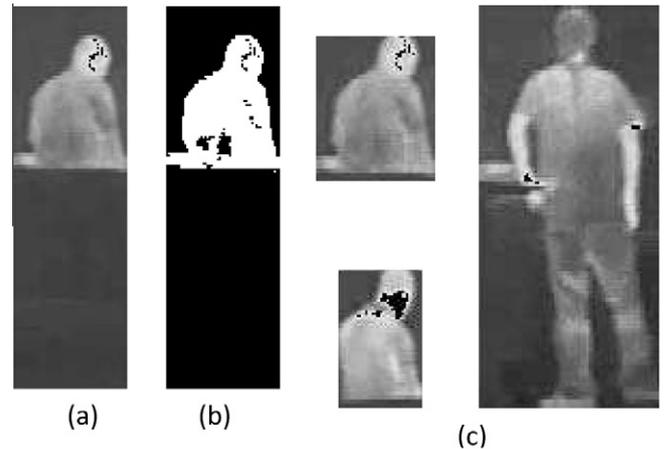


Fig. 3. (a) Input sub-ROI, (b) Binarized sub-ROI, (c) Rows adjustment to delimit three human candidates.

Table 2
Refined people candidates blobs list.

κ	x_{start}	y_{start}	x_{end}	y_{end}	Area
1	339	286	395	354	3808
2	396	270	481	458	15980
3	298	289	338	354	2600

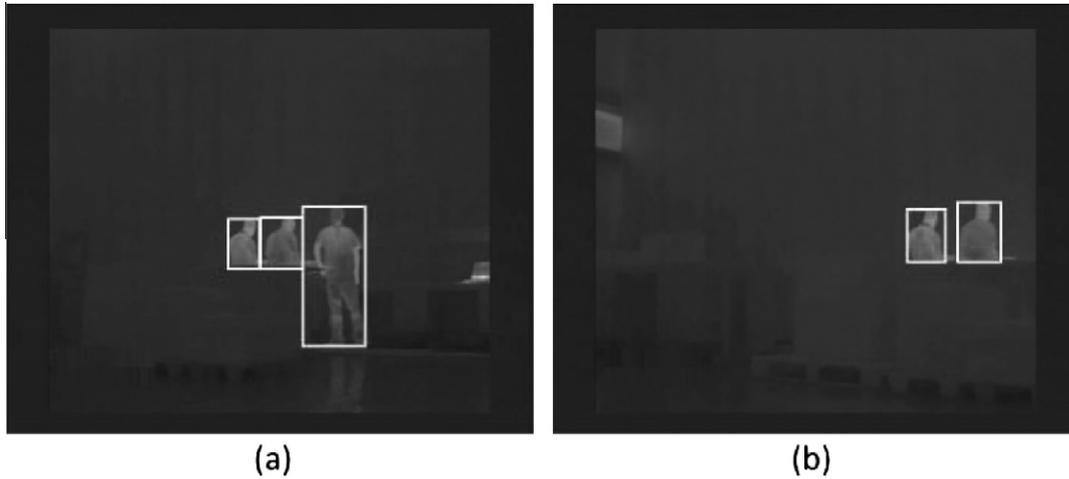


Fig. 4. Some indoor output images.

height of the human contained. Rows adjustment is performed for each new sub-ROI, $sR_{i,\alpha}$, generated by the previous columns adjustment, by applying a new threshold, θ_h .

The calculation is applied separately on each sub-ROI to avoid the influence of the rest of the image on the result. This threshold takes the value of the sub-ROI average gray level, $\theta_h = \overline{sR_{i,\alpha}}$. Thus, sub-ROI $sR_{b,i,\alpha}$ is binarized in order to delimit its upper and lower limits, obtaining $sR_{i,\alpha}$, as shown in Eq. (12) similar to Eq. (5).

$$sR_{b,i,\alpha}(i,j) = \begin{cases} \min, & \text{if } sR_{i,\alpha}(i,j) \leq \theta_h \\ \max, & \text{otherwise} \end{cases} \quad (12)$$

After this, a closing is performed to unite spots isolated in the binarization, getting $sR_{c,i,\alpha}$ (see Fig. 3(b)).

$$sR_{c,i,\alpha}(i,j) = sR_{b,i,\alpha}(i,j) \bullet \begin{vmatrix} 0 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{vmatrix} \quad (13)$$

Next, $sR_{c,i,\alpha}$ is scanned, searching pixels with values superior to \min . The upper and lower rows of the human are equal to the first and last rows, respectively, containing pixels with a value set to \max . The final result, assigned to new ROIs, \mathfrak{R}_k , may be observed in Fig. 3(c). The blobs associated to the split ROIs are enlisted into the original blobs list, L_B (see Table 2).

3. Data and results

The people segmentation algorithm proposed has been tested on a series of indoor video sequences captured by a FLIR camera. The aim has been to extensively show the robustness of the proposal in only obtaining humans through the segmentation of warmer spots and the refinement to eliminate non-human objects. Fig. 4 shows a couple of output images of the infrared human detection algorithm. As you may easily observe, in both frames presented, humans are perfectly segmented.

Fig. 4(a) shows the clear separation of three very close people in the running example used in the previous section. Notice that two of the three humans are partially occluded, and, nonetheless, they are correctly segmented. It is also worthwhile to highlight that the shadow of the non-occluded human is eliminated from the human body. In the same figure, a laptop (at the right of the people) is not segmented as a human, although there is sufficient heat emitted by the computer to confuse the algorithm. Also, consider Fig. 4(b),

where a window at the left of the image has not provided a false positive.

The performance results in terms of real-time capability of the algorithms described are also excellent, as the method works with 6 frames per second in the worst situation, which corresponds to the processing of large typical 768×576 pixel images provided by a FLIR camera (see Fig. 5). When processing smaller images are processed, the performance reaches 13 frames per second for 537×403 pixel images and 26 frames per second for 384×288 pixel images in the test performed.

Of course, we have also tested our proposal with a well-known indoor infrared video benchmark, namely the “Indoor Hallway Motion” sequence included in dataset 5 “Terravic Motion IR Database” provided within the OTCBVS Benchmark Dataset Collection.

Firstly, in Fig. 6 the results of applying the algorithm to a sequence where a person is crossing the scene from right to left. Fig. 6(a) shows the scenario. Notice the presence of a hot spot in the image that could lead to confusion, but the algorithm does not make any mistake. Fig. 6(b) shows the first frame where the person is detected. The detection works fine – a hit of a 100% – through all the frames (see Fig. 6(c)) until the person reaches the hot spot region. Here (Fig. 6(d)), the human region also covers the hot spot. From this moment on, and until the person exits the scene, the performance of the person segmentation hit is a 98% (e.g. Fig. 6(e) and (f)).

The second example shows another person crossing from right to left, while, at the same time, a second human is slowly

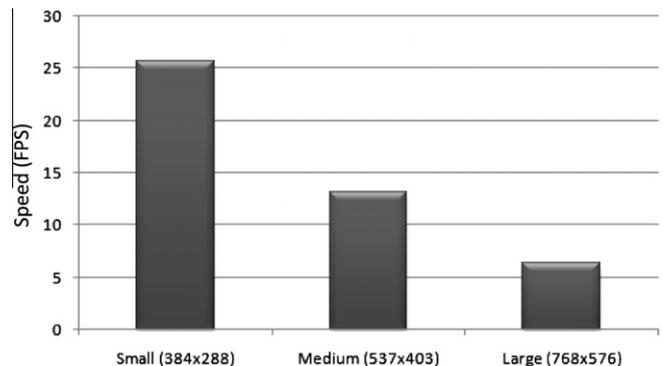


Fig. 5. People segmentation algorithm (speed in fps).

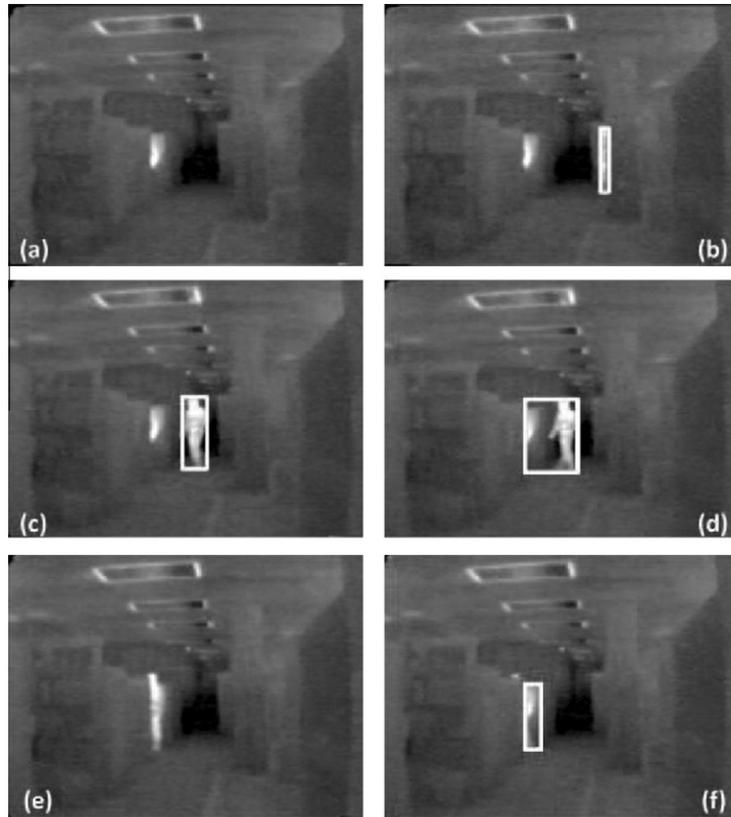


Fig. 6. Crossing from right to left. (a) Frame #220, (b) Frame #250, (c) Frame #290, (d) Frame #292, (e) Frame #324, (f) Frame #333.

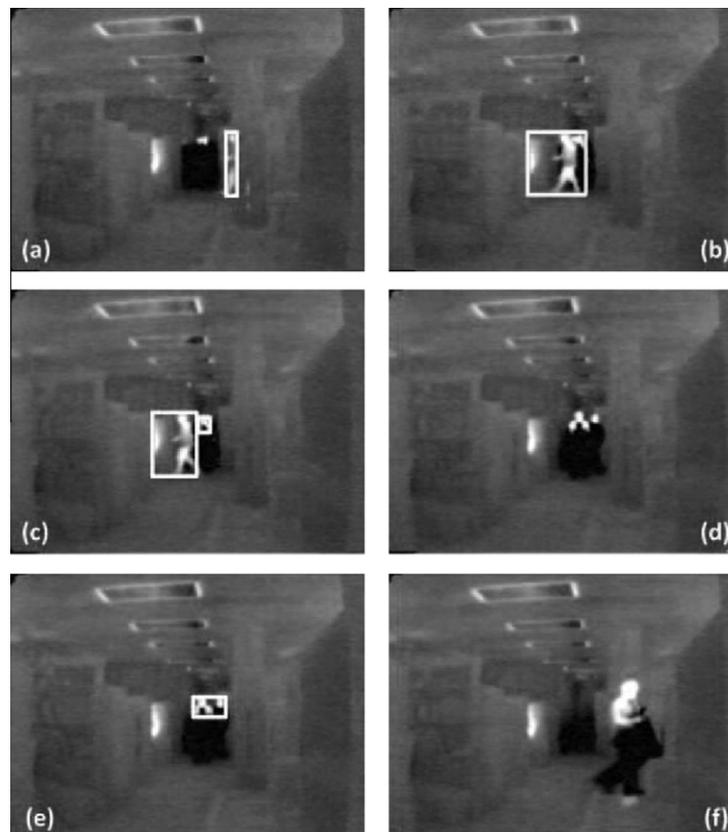


Fig. 7. Crossing from right to left, and approaching from rear to front. (a) Frame #2326, (b) Frame #2357, (c) Frame #2362, (d) Frame #2420, (e) Frame #2411, (f) Frame #2530.

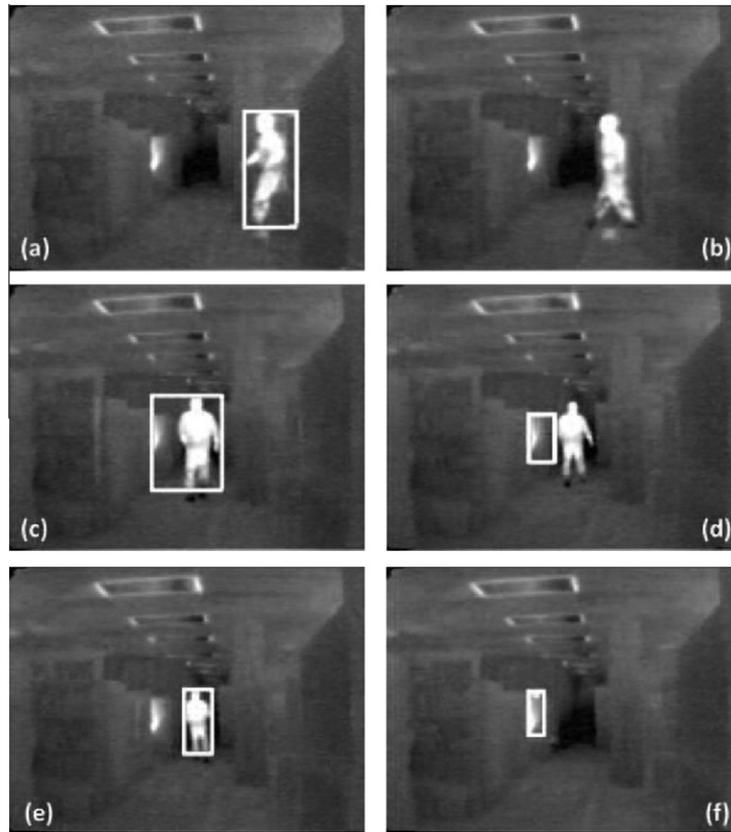


Fig. 8. Walking from front right to rear left. (a) Frame #5135, (b) Frame #5150, (c) Frame #5175, (d) Frame #5209, (e) Frame #5229, (f) Frame #5264.

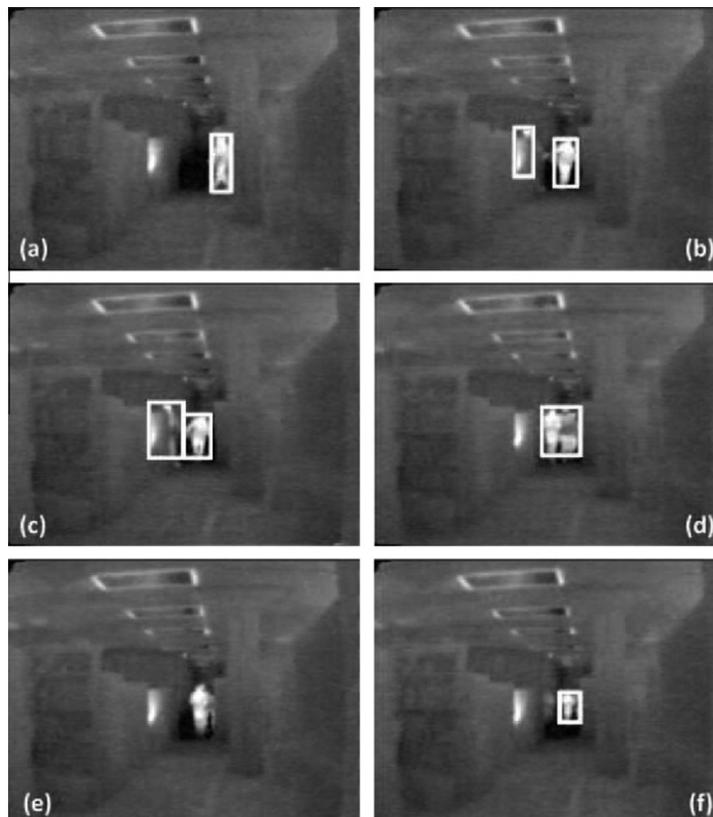


Fig. 9. Two persons walking to the rear. (a) Frame #8268, (b) Frame #8293, (c) Frame #8297, (d) Frame #8344, (e) Frame #8377, (f) Frame #8478.

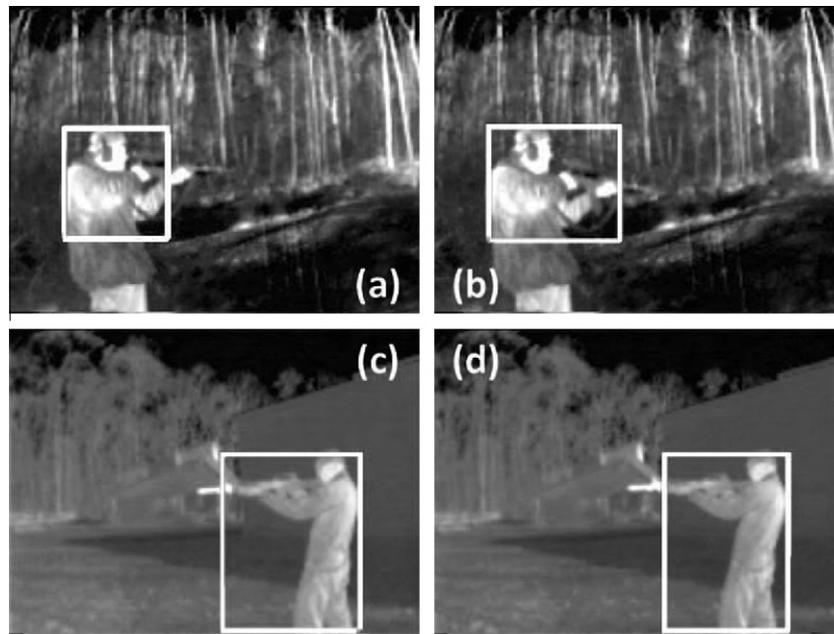


Fig. 10. Some results on outdoor infrared images.

approaching the camera from the rear. As you may observe (Fig. 7(a)–(c)), the first human is perfectly segmented. The person approaching is wearing clothes that do not emit any heat. Thus, only the head is visible, and it is very difficult to detect such a spot as belonging to a human. Nonetheless, in a 63% of the frames where the partially occluded human is present, he is detected (see Fig. 7(d) and (e)). Unfortunately, when this person is too close to the camera, it is not detected as a human (see Fig. 7(f)).

The third example shows a human entering the scene from the front right and walking to the rear left (see Fig. 8). Now, this person shows an abnormal and uniform high heat level. This is a challenge for our proposal, as we try to locate the human head in relation with the body (in terms of heat and size). The performance of the algorithm in this sequence is about a 60% of hits when the person is close to the camera (see Fig. 8(a)–(c)) and grows to a 90% when the human goes far, as shown in Fig. 8(e) and (f). Fig. 8(d) shows one of the three frames where a false positive is gotten in a total of 140 frames of the sequence.

Finally, we show in Fig. 9(a) sequence of two people entering the scene from different sides, and walking together to the rear. As shown in Fig. 9(a)–(c), the segmentation process is efficient. From frame #8344 on, corresponding to Fig. 9(d), the humans are not correctly divided, as they are too close to the rear. As shown in Fig. 9(e) and (f), until the humans disappear from the scene, the segmentation of the group is right for a 47%.

Lastly, for the purpose of discovering how good our algorithms could work in outdoor scenarios, we tested our algorithm on some outdoor benchmarks (also from the OTCBVS Benchmark Dataset Collection – “Weapon Presence Detection” and “Weapon Discharge Detection”, IR05 and IR03, respectively, of dataset 6 “Terravic Weapon IR Database”). Just by modifying a few parameters of the algorithm, we came to the conclusion that the proposal is very encouraging, as shown by the results offered in Fig. 10 and Table 3.

Table 3
Outdoor human detection results.

Benchmark	Frames	Hits	False positives	False negatives
IRG03	480	98%	0%	2%
IRG05	297	100%	0%	0%

4. Conclusions

A new approach to real-time people segmentation through processing images captured by an infrared camera has been extensively described. The proposed algorithm starts detecting human candidate blobs processed through traditional image thresholding techniques. Afterwards, the blobs are refined with the objective of solving the question if each blob contains one single human candidate or has to be divided into smaller blobs. If the blob contains more than one possible human, the blob is divided to fit each new candidate in height and width.

The results obtained so far in indoor and (initial) outdoor scenarios are promising. We have been able of testing our person segmentation algorithms on our proper testbeds, as well as on very well-known datasets.

Acknowledgements

This work was partially supported by the Spanish Ministerio de Ciencia e Innovación under project TIN2007-67586-C02-02, and by the Spanish Junta de Comunidades de Castilla-La Mancha under projects PII2I09-0069-0994, PII2I09-0071-3947 and PEII09-0054-9581.

References

- Benet, G., Blanes, F., Simó, J. E., & Pérez, P. (2002). Using infrared sensors for distance measurement in mobile robots. *Robotics and Autonomous Systems*, 40(4), 255–266.
- Bhanu, B., & Han, J. (2002). Kinematic-based human motion analysis in infrared sequences. In *Proceedings of the sixth IEEE workshop on applications of computer vision* (pp. 208–212).
- Davis, J. W., & Sharma, V. (2007). Background-subtraction in thermal imagery using contour saliency. *International Journal of Computer Vision*, 71(2), 161–181.
- Fernández-Caballero, A., Gómez, F. J., & López-López, J. (2008). Road-traffic monitoring by knowledge-driven static and dynamic image analysis. *Expert Systems with Applications*, 35(3), 701–719.
- Fernández-Caballero, A., López, M. T., Mira, J., Delgado, A. E., López-Valles, J. M., & Fernández, M. A. (2007). Modelling the Stereovision-correspondence – Analysis task by lateral inhibition in accumulative computation problem-solving method. *Expert Systems with Applications*, 33(4), 955–967.
- Fernández-Caballero, A., López, M. T., & Saiz-Valverde, S. (2008). Dynamic stereoscopic selective visual attention (DSSVA): Integrating motion and shape

- with depth in video segmentation. *Expert Systems with Applications*, 34(2), 1394–1402.
- García, M. A., & Solanas, A. (2004). Estimation of distance to planar surfaces and type of material with infrared sensors. In *Proceedings of the 17th international conference on pattern recognition* (Vol. 1, pp. 745–748).
- Iwasawa, S., Ebihara, K., Ohya, J., & Morishima, S. (1997). Realtime estimation of human body posture from monocular thermal images. In *Proceedings of the 1997 IEEE computer society conference on computer vision and pattern recognition* (pp. 15–20).
- Jung, S. -H., Eledath, J., Johansson, S., & Mathevon, V. (2007). Egomotion estimation in monocular infra-red image sequence for night vision applications. In *IEEE workshop on applications of computer vision* (p. 8).
- López, M. T., Fernández-Caballero, A., Fernández, M. A., Mira, J., & Delgado, A. E. (2006). Visual surveillance by dynamic visual attention method. *Pattern Recognition*, 39(11), 2194–2211.
- López, M. T., Fernández-Caballero, A., Mira, J., Delgado, A. E., & Fernández, M. A. (2006). Algorithmic lateral inhibition method in dynamic and selective visual attention task: Application to moving objects detection and labelling. *Expert Systems with Applications*, 31(3), 570–594.
- Mira, J., Delgado, A. E., Fernández-Caballero, A., & Fernández, M. A. (2004). Knowledge modelling for the motion detection task: The algorithmic lateral inhibition method. *Expert Systems with Applications*, 27(2), 169–185.
- Nanda, H., & Davis, L. (2002). Probabilistic template based pedestrian detection in infrared videos. In *Proceedings of the IEEE intelligent vehicle symposium* (Vol. 1, pp. 15–20).
- Pavón, J., Gómez-Sanz, J., Fernández-Caballero, A., & Valencia-Jiménez, J. J. (2007). Development of intelligent multi-sensor surveillance systems with agents. *Robotics and Autonomous Systems*, 55(12), 892–903.
- Xu, F., Liu, X., & Fujimura, K. (2005). Pedestrian detection and tracking with night vision. *IEEE Transactions on Intelligent Transportation Systems*, 6(1), 63–71.
- Yilmaz, A., Shafique, K., & Shah, M. (2003). Target tracking in airborne forward looking infrared imagery. *Image and Vision Computing*, 21(7), 623–635.