# Automatic People Detection and Counting for Athletic Videos Classification

C. Panagiotakis[1]    E. Ramasso[2]    G. Tziritas[1]    M. Rombaut[2]    D. Pellerin[2]

[1] Computer Science Department
University of Crete
Heraklion, Greece

[2] GIPSA-lab
Images and Signal Department
Grenoble, France

## Abstract

We propose a general framework that focuses on automatic individual/multiple people motion-shape analysis and on suitable features extraction that can be used on action/activity recognition problems under real, dynamical and unconstrained environments. We have considered various athletic videos from a single uncalibrated, possibly moving camera in order to evaluate the robustness of the proposed method. We have used an easily expanded hierarchical scheme in order to classify them to videos of individual and team sports. Robust, adaptive and independent from the camera motion, the proposed features are combined within Transferable Belief Model (TBM) framework providing a two level (frames and shot) video categorization. The experimental results of $97\%$ individual/team sport categorization accuracy, using a dataset of more than 250 videos of athletic meetings indicate the good performance of the proposed scheme.

## 1 Introduction

Automatic indexing based on human motion analysis and action/activity recognition under real, dynamical and unconstrained environments is key of importance because can be applied in many areas such as database management, surveillance or human-computer interface. Human motion analysis consists in [1] *detection, tracking* and *recognition*. Group detection and/or counting is generally embedded in the *detection and tracking* processes. Object detection and tracking in complicated environments is still the key problem of the visual surveillance and it is becoming an important issue in several applications such as camera based surveillance and human machine interaction. The detection and tracking algorithms are challenged by occluding and fast/complicated moving objects, as well as illumination changes.

Concerning the 2-D approaches, Wang et al. [2] propose a method to recognize and track a walker using 2D human model and both static and dynamic cues of body biomet-

rics. Moreover, many systems use Shape-From-Silhouette methods to detect and track the human in 2D [3] or 3D space [4]. The silhouettes are easy to extract providing valuable information about the position and shape of the person. When the camera is static, background subtraction techniques can give high accuracy measures of human silhouettes by modeling and updating the background image [5]. Otherwise, when the camera is moving, camera motion estimation methods [6] can locate the independently moving objects. The system called W4 [5] is based on a statistical-background model to locate people and their parts (head, hands, feet, torso, etc.) using static cameras and allowing multiple person groups. Rabaud and Belongie [7] present a method for counting moving objects without tracking them based on a highly parallelized version of the KLT tracker. It is performed in crowding situations where the tracking does not make sense to perform. Figueroa et al. [8] propose a system of tracking soccer players using multiple static cameras. The occlusions have been treated by splitting segmented blobs based on morphological operators and a backward and forward graph representation based on human shape, motion and color features. However, in a real soccer game, there are crowd situations, where the people should be manually tracked.

Most of aforementioned schemes are semi-automatic and assume static camera and constrained indoor environments, analyzing high quality silhouettes (Figure 1(b)). However, the estimated silhouettes from real and unconstrained environments are in low quality (Figures 1(d) and 1(f)). The goal of this research is to propose a general hierarchical scheme that can be used on action/activity recognition problems under real, dynamical and unconstrained environments, based on low level shape features. In order to evaluate the robustness of the proposed scheme, we have applied it on various athletic videos from a single uncalibrated and possibly moving camera. The automatic analysis of these videos is a challenging problem due to the complex and fast motions of the athletes and to the unconstrained changes in the environment of athletic meetings. In previous work, a novel architecture utterly based on Transfer-
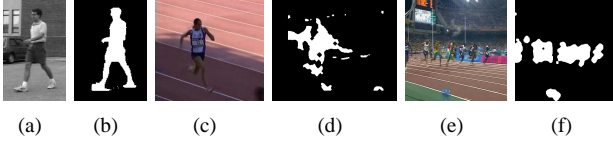
Figure 1: **(a), (b)** Original image and the silhouette estimated by the method of [5] under static camera. **(c), (e)** Original images and **(d), (f)** the corresponding silhouettes estimated by the method of [10] under moving camera.
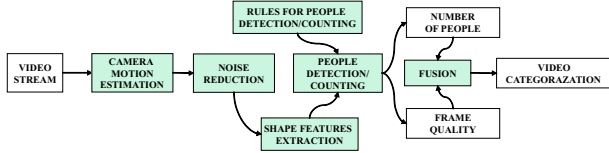


Figure 2: Schema of the proposed system architecture.

able Belief Model [9], was proposed [3, 10] for individual motion analysis and action/activity recognition in athletic sports videos.

The goal of the proposed method is to classify a video into individual (individual sport such as high jump and long jump) and multiple people (team sport such as running and hurdling). As well, the system detects and counts the number of people in videos. A reliability factor is estimated at each frame in order to quantify the quality of the classification which is taken into account for decision concerning the number of people. In this framework, no initialization step is required and no assumption or knowledge is assumed about the number of people and their motion in the scene. The proposed framework can be decomposed into several main modules illustrated in Figure 2. First, silhouettes are computed using a camera motion estimation method [6] and shape based features are extracted (Section 2). Next, people detection and counting are performed (Section 3). Finally, the video is classified to individual sport and team sport video, based on a TBM fusion process (Section 4). Experimental results are given in Section 5. Conclusions and discussion are provided in Section 6.

## 2. Features Extraction

The binary silhouettes, estimated by the camera motion estimation method [6][1], probably contain objects that do not follow the athletes motion, e.g. fake objects in the background (see Figure 1). We assume that the camera tries to track the humans (athletes), which is a tenable assumption

---

[1]An affine model is used to describe the camera motion. The above method, that we use, was implemented by the Vista Team of IRISA and has the advantage to take into account the global variation of illumination thus it is adapted for real videos.
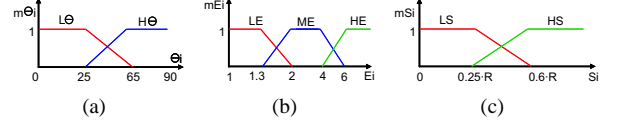


Figure 3: From numerical features to belief. **(a)** Angle, **(b)** Eccentricity, **(c)** Area.

because the athlete is the object of interest, so the athletes are about in the same position in a short time window. If we suppose that the noise is appearing in random positions (white noise) over the time, then a lowpass time filter can remove the noise (see gray pixels of Figure 4).

Shape features are computed in order to detect both humans and groups and to count people in groups. We compute for each object $O_t^i$, its major axis angle $\theta_i$, its eccentricity $\varepsilon_i$ and its normalized area $s_i$. These features are uncorrelated, independent from camera view, and their values can be estimated robustly under low quality silhouettes.

The angle of the object $O_t^i$ major axis $\theta_i$ is defined by the three second order moments. This angle shows the main orientation of the object. In the whole paper, angles are measured in degrees. The robustness of $\theta_i$ estimation is determined by the object's eccentricity. The eccentricity ($\varepsilon_i \geq 1$) is defined by the ratio between the two principal axes of the best fitting ellipse, measuring how thin and long a region is. If $\varepsilon_i$ is close to one, then $\theta_i$ will be unspecified. The area feature $s_i$ should be normalized in order to be independent from both image size and distance of the object from the camera. Generally, it holds that the area of interest concerns the athletes, that are tracked by the camera. These athletes normally have similar distances from the camera. Therefore, $s_i$ is defined as ratio between an object area ($|O_t^i|$) and the mean object area $\frac{\sum_k |O_t^k|}{N_t}$, where $N_t$ denotes the estimated number of people at frame $t$ (Section 4.1).

The three proposed features ($\theta_i$, $\varepsilon_i$ and $s_i$) are simple and well understandable. Thus, they can be easily converted into beliefs using a fuzzy-sets inspired symbolic representation (Figure 3). Appropriate table of rules can then be used for people detection and counting as presented in Tables 1-2 (and discussed further). The proposed numeric-to-symbolic conversion is presented in Figure 3, where $L$ is used for low value, $M$ for medium values and $H$ for high values. Figure 3(a) presents the angle $\Theta_i$ numeric-to-symbolic conversion, with $\Theta_i = min(\theta_i, 180 - \theta_i)$, $\Theta_i \in [0, 90]$. Figure 3(b) presents the numeric-to-symbolic conversion of eccentricity. The red, blue and green curves correspond to the probability of $LE$, $ME$ and $HE$ respectively. Figure 3(c) presents the numeric-to-symbolic conversion of area feature $s_i$. Two beliefs are concerned: low area ($LS$), which is true for $s_i \leq 0.25 \cdot R^2$ indicating a little area objects (possibly

---

[2]$R$ is an adaptive factor, denoting the probability of an object, which

noise), and high area, indicating an object of normal area (possibly humans).

# 3. People Detection and Counting

## 3.1. People detection

This method removes noise objects, that can not be removed by the noise reduction procedure, using their shape features. We have used the rules of Table 1 in order to detect and remove such objects, combining the symbolic beliefs. This table can be estimated by a learning stage using an EM procedure for instance. Using this table, the probability of human $P_r(O_t^i\ is\ H)$ can be estimated by the cells where the probability of $H$ is positive. An object $O_t^i$ will be detected as human if $P_r(O_t^i\ is\ H) > 0.5$. Otherwise, it will be detected as noise.

|  | LE | ME | HE |
|---|---|---|---|
| LS, LΘ | $N$ | $N$ | $N$ |
| LS, HΘ | $N$ | 0.15/$H$ 0.85/$H$ | $N$ |
| HS, LΘ | 0.2/$N$ 0.8/$H$ | $H$ | $H$ |
| HS, LΘ | 0.1/$N$ 0.9/$H$ | $H$ | 0.05/$N$ 0.95/$H$ |

Table 1: The table rules for human/noise detection, $N$, $H$ denote *noise*, *humans*, respectively.

The notation $0.2/N$ in Table 1 means that the $20\%$ of the belief is assigned to proposition $N$.

## 3.2. People counting

The people counting procedure is executed for each human (blob) detected object and is based on the assumption that each human major axis (in the most time) is mainly vertical. Thus, an individual object probably has high angle and medium eccentricity. Otherwise, the object is probably a group of people containing two, three or more people. Using rules of Table 2, where the proposed features are combined, the number of people per object can be estimated, where $K$ denotes the number of people (real value) in groups estimated by an algorithm described hereafter.

The number of people $K_i$ (real value) of a horizontally directed object $O_t^i$ is estimated by using its eccentricity ($\varepsilon_i$) and its area ($|O_t^i|$). According to the objects' surfaces, the most possible value of $K_i$ is $s_i$. The object is horizontally directed, so we can use the group model of Figure 5(a) in order to compute the mean eccentricity per human $e_h = \frac{H}{L_h}$, where $L = \sqrt{|O_t^i| \cdot \varepsilon_i}$ (group length), $H = \frac{|O_t^i|}{L}$

---
has $s_i > 0.05$, to be a human object.

(group height) and $L_h = \frac{L}{s_i}$ (human length). If $e_h$ is higher than four, which is the maximum individual eccentricity, the number of humans will be recomputed by enforcing the eccentricity per human to be four. Thus, it holds that, $L_h = \frac{H}{4}$ and $K_i = \frac{L}{L_h}$. Figure 4 illustrates the results of people detection and counting algorithm. The little black boxes corresponds to the mass centers of the detected humans.
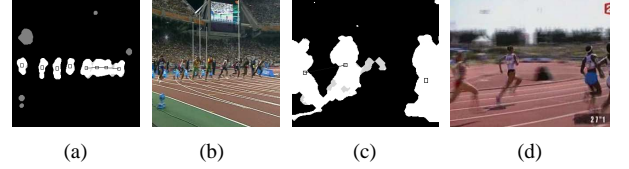

(a)    (b)    (c)    (d)

Figure 4: **(a)** Four individual and one group of four people are detected. **(b)** The original image of **(a)** is shown. **(c)** An individual and a group of two people are detected. **(d)** The original image of **(c)** is shown.

## 3.3. Quality factor estimation

A measurement of frame quality (reliability) factor $Q_t$ can be estimated using the probability of the decisions (human/noise decision and counting decision) results. If the decisions are taken with low probabilities then $Q_t$ should be low, otherwise $Q_t$ should be high. Let $P_r^{HN}(O_t^i)$ denotes the decision probability of the object $O_t^i$ to be human or noise. Let $P_r^{NP}(O_t^i)$ be the decision probability concerning the number of people in the object $O_t^i$. $Q_t \in [0,1]$ is estimated by the product of the expected values ($E_i$) of $P_r^{HN}(O_t^i)$, $P_r^{NP}(O_t^i)$ over the objects:

$$Q_t = E_i(P_r^{HN}(O_t^i)) \cdot E_i\left(\frac{P_r^{NP}(O_t^i)}{\sqrt{\max(K_i,1)}}\right) \qquad (1)$$

$P_r^{NP}(O_t^i)$ is divided with the square root of the number of detected people $\sqrt{\max(K_i,1)}$ in object $O_t^i$, because the accuracy of people counting procedure decreases, as the number of people increases (occlusions are appeared). The use of $\sqrt{\max(K_i,1)}$ improves slightly the categorization accuracy. $Q_t$ will be used on video categorization scheme. Figure 5(c) presents $Q_t$ numeric-to-symbolic conversion. There are three beliefs for quality factor: bad quality (**Bad**), unknown quality (**Bad** $\cup$ **Good**) and high quality (**Good**).

# 4. Video Categorization Scheme

The results of people counting procedure and the frame quality factor are fused using TBM framework in order to discriminate the video of individual sport ($I$) and team sport ($T$). We have used the TBM framework, since it is more general than probabilities and explicitly defines the conflict

|  | **LE** | **ME, HE** |
|---|---|---|
| **HS, LΘ** | $1, 0.3 + \max(1 - \|K - \min(1, K)\|, 0)$ <br> $2, 1.0 + \max(1 - \|K - 2\|, 0)$ <br> $3, 0.8 + \max(1 - \|K - \max(3, K)\|, 0)$ | $\lfloor K+1 \rfloor, 1 - K + \lfloor K+1 \rfloor$ <br> $\lfloor K \rfloor$, remainder |
| **HS, HΘ** | $1, 0.6 + \max(1 - \|K - \min(1, K)\|, 0)$ <br> $2, 1.0 + \max(1 - \|K - 2\|, 0)$ <br> $3, 0.8 + \max(1 - \|K - \max(3, K)\|, 0)$ | $1, 0.6 + 0.4 * \max(1 - \|K - 1\|, 0)$ <br> $2$, remainder |

Table 2: The table rules for people counting.

and doubt. The classification concerns two classes: video of individual sport ($I$) and team sport ($T$). Therefore, $\Omega = \{I, T\}$ is the frame of discernment of the classification. A basic belief assignment (BBA) [9] $m_t^\Omega$ at frame $t$ is defined on the set of propositions $2^\Omega = \{\emptyset, I, T, I \cup T\}$, where $\emptyset$ and $I \cup T$ correspond to the conflict and doubt respectively. $m_t^\Omega : 2^\Omega \to [0, 1], X \to m_t^\Omega(X)$ and by construction it holds that $m_t^\Omega(\emptyset) = 0$, and $\sum_{X \subseteq \Omega} m_t^\Omega(X) = 1$. A value $m_t^\Omega(X)$ is a basic belief mass which expresses a confidence proposition $X \subseteq \Omega$ according to a given feature but does not imply any additional claims regarding subsets of $X$. It is the fundamental difference with probability theory.

## 4.1. Number of people estimation

The number of people $N_t$ (real value) at frame $t$ is robustly estimated using quality factor:

$$N_t = (1 - Q_t) \cdot N_{t-1} + Q_t \cdot TP_t \tag{2}$$

where $TP_t$ (integer value) denotes the sum of the number of the detected people (for all blobs) at frame $t$ using the rules of Table 2 (see Section 3.2). The estimation of number of people by $N_t$ is more robust than by $TP_t$, since it takes into account the quality factor. Figure 5(b) presents $N_t$ numeric-to-symbolic conversion which provides the BBA $m_{N_t}$ defined for three sets: low number of people ($I$), middle number of people ($I \cup M$), and high number of people ($M$).
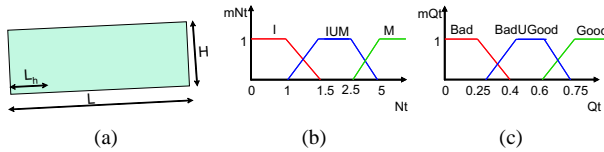


Figure 5: **(a)** Group model. **(b)**, **(c)** From numerical features to belief. **(b)** Number of people, **(c)** Quality factor.

Beliefs concerning both the number of people and quality (obtained from the previous steps), i.e. $m_{N_t}$ and $m_{Q_t}$, are combined using the rules in Table 3 resulting in the belief mass denoted $m_t^\Omega$.

|  | **I** | **I ∪ M** | **M** |
|---|---|---|---|
| **Bad** | $I \cup T$ | $I \cup T$ | $I \cup T$ |
| **Bad ∪ Good** | $0.5/I$ <br> $0.5/I \cup T$ | $I \cup T$ | $0.5/T$ <br> $0.5/I \cup T$ |
| **Good** | $I$ | $I \cup T$ | $T$ |

Table 3: Table of rules for individuals/groups detection.

## 4.2. Short time belief and decision

In order to take temporal aspects of belief $m_t^\Omega$ into account, we combine it with the previous belief of the system concerning the category of the video using the TBM conjunctive rule of combination [9]. By doing it, we assume that the belief does not evolve between two successive frames. Note that a method was proposed in [11] for belief functions filtering and that includes a model of evolution between successive frames. The fusion process is performed frame by frame for each proposition $X$ yielding a new local mass $\hat{m}_t^\Omega(X)$:

$$
\begin{aligned}
\hat{m}_t^\Omega(X) &= \hat{m}_{t-1}^\Omega \textcircled{$\cap$} m_t^\Omega(X) \\
&= \sum_{C \cap D = X} \hat{m}_{t-1}^\Omega(C) \cdot m_t^\Omega(D)
\end{aligned}
\tag{3}
$$

This belief mass quantifies the system's belief concerning the categorization of the video in the frame of discernment $\Omega = \{I, T\}$ previously defined. Using the aforementioned fusion process, the mass on the empty set ($\hat{m}_t^\Omega(\emptyset)$), called conflict, is going to increase to one, while the masses on the other propositions are going to decrease to zero. This effect is due to the fact that the empty set is abortive by the $\textcircled{$\cap$}$-rule. When the conflict is high, the trapezes used in the numeric to symbolic conversion are modified manually in order to decrease the conflict (by adding doubt for instance). When the conflict is not too high, we have used the Dubois & Prade's conflict redistribution rule [12] in order to manage conflict yielding to: $\hat{m}_t^\Omega(\emptyset) = 0$:

$$\hat{m}_t^\Omega(C \cup D) = \sum_{C \cap D = \emptyset} \hat{m}_{t-1}^\Omega(C) \cdot m_t^\Omega(D) \tag{4}$$

This rule is called adaptive because it redistributes the value of conflict onto the union of hypotheses that cause conflict.

## 4.3. Final decision

The final decision concerning the whole video sequence is taken by "equivalent" fusion of the beliefs at each frame. Thereby, at frame $t$, the mean mass $\bar{m}_t^\Omega(X)$ of the proposition $X$ is computed by getting the mean of the local decision mass $\hat{m}_k^\Omega(X)$ over the frames $\{1, 2, \cdots, t\}$:

$$\bar{m}_t^\Omega(X) = \frac{1}{t} \cdot \sum_{k=1}^{t} \hat{m}_k^\Omega(X) \tag{5}$$

Finally, the decision is taken using the pignistic probability (BetP) proposed by Ph. Smets [13]. The above decision rule is equivalent with the selection of the proposition $X$ with the highest mean mass $\bar{m}_t^\Omega(X)$ (because the frame of discernment is made of two hypotheses). Using the aforementioned scheme, the value of the selected mean mass provides a final decision confidence value. This value corresponds to the probability of the final decision.

## 5. Experimental Results

We have tested the proposed algorithm on a data set containing 252 athletic videos captured from broadcast TV: 161 video sequences from individual sports like pole vault, high jump, shot, javelin, etc. and 91 video sequences from team sports like running and hurdling. The database is characterized by its heterogeneity with a panel of view angles as well as unconstrained indoor or outdoor environments (other moving people can appeared) and athletes (male, female with different skills and skin colors).

The accuracy of the team sports detection was $96.9\%$ (156/161) and the accuracy of the individual sports was $96.7\%$ (88/91). We have performed several tests in order to make comparisons between the proposed scheme versus several variations of this scheme. First, we tested the proposed scheme without using quality criterion (setting quality factor equal to one), getting $93.8\%$ for individual sports and $97.8\%$ for team sports. We tested the proposed scheme without using TBM framework, deciding using a threshold $(Th_r)$ on the mean of $N_t$ over the frames, getting about $3\%$ less performance. Conclusively, the aforementioned comparisons show the importance of using quality $Q_t$, TBM framework and the robustness of the proposed features under several decisions rules. Some videos and experimental results by the proposed method are available at the Web addresses `www.csd.uoc.gr/~cpanag/ DEMOS/actionActivityRecognition.htm` and `www. lis.inpg.fr/pages_perso/ramasso`.

Figures 6 and 7 show frames from the original sequences and the corresponding results of the proposed scheme. The people in group are connected with straight lines. In Figure 6, two athletes are initially appearing in the scene making the method confusing to decide (at first frames). Finally, the camera tracks one of them and the system responds that it was an individual sport video ($N = 1$ at the end (Fig. 6-xiii), $N$ varies between 1 and 2.5 within the sequence). The belief mass $\hat{m}_t^\Omega(X)$ gives an instant decision for a current frame. According to the $\hat{m}_t^\Omega(X)$ until the frame 35 multiple people are appearing in the scene, which is very close to the ground truth. The global decision for a period can be taken using $\bar{m}_t^\Omega(X)$. According to this mass, after the 70th frame, the video is classified as an individual sport, since it contains more frames of single athlete rather than multiple athletes. Figure 7 illustrates a 100 m running video, which is correctly classified into team sports. First, the 8 athletes are separated providing high accuracy results to the people counting procedure and high values on $Q_t$. After the middle of the sequence, a lot of occlusions and bad quality silhouettes are appearing. The occluded athletes correspond to one or two groups of people, and at the same time $Q_t$ has low values. This example shows the accuracy of people counting procedure under several conditions ($N = 8$ at the beginning of the sequence (Fig. 7-xiii), then $N$ decreases due to the view angle). It show the usefulness of $Q_t$ in order to be able to give a confidence value about the people counting at each frame. Thus, the value of quality is very close to what a human expert will decide for a quality function, since it is maximized when athletes are well detected (without noise). Accurate assessment of number of people is hard to perform because one needs to build a ground truth and, besides, single view camera used in our application limits occlusion handling (see Fig. 7). However, the evolution of $N_t$ on both figures shows that the system can detect precisely the number of people.

## 6. Conclusion

We have proposed a shape based method for unsupervised-automatic people detection and counting applied to athletic videos in order to classify them to videos of individual sports and team sports. Robust, adaptive, independent from the camera motion and well understandable by humans features, are estimated using silhouettes. Finally, the features are combined within Transferable Belief Model (TBM) framework for video categorization yielding at the same time confidence values about the final decision.

The main contribution of this work concerns the definition of appropriate robust features and the TBM based fusion of them, using a quality function, yielding high performance results without any given feature or initialization under low quality - real conditions videos. The proposed method can be easily applied to other types of videos. In particular, the trapezes as well as tables of rules could be
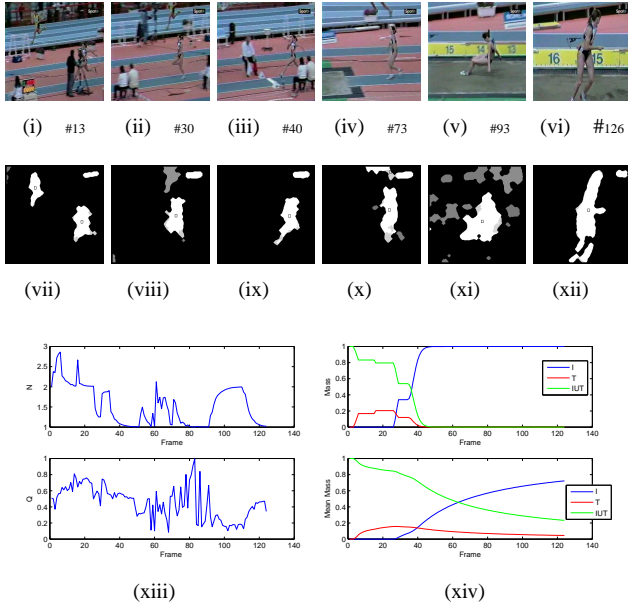
Figure 6: **(i)**, $\cdots$, **(viii)** Triple jump original sequence with 126 frames. **(ix)**, $\cdots$, **(xvi)** The results of the people detection and counting procedure. The small black boxes correspond to the mass center detected humans. **(xvii)** $N_t$, $Q_t$. **(xviii)** The belief masses $\hat{m}_t^\Omega(X)$, $\bar{m}_t^\Omega(X)$.



Figure 7: **(i)**, $\cdots$, **(viii)** Original running sequence with 243 frames. **(ix)**, $\cdots$, **(xvi)** The results of the people detection and counting procedure. Small black boxes: mass center of detected humans. A group of people is detected when the boxes are connected with a line. **(xvii)** $N_t$, $Q_t$. **(xviii)** The belief masses $\hat{m}_t^\Omega(X)$, $\bar{m}_t^\Omega(X)$.

estimated according to the type of videos using for instance majority rule algorithm [14] but this need an heavy step for manual annotation in order to prepare the learning set.

# Acknowledgments

# References

[1] L. Wang, W. Hu, and T. Tan, "Recent developments in human motion analysis," *Pattern Recognition*, vol. 36, no. 3, pp. 585–601, 2003.

[2] L. Wang, H. Ning, and T. Tan, "Fusion of static and dynamic body biometrics for gait recognition," *IEEE Trans. on CSVT*, vol. 14, no. 2, pp. 149–158, 2004.

[3] C. Panagiotakis, E. Ramasso, G. Tziritas, M. Rombaut, and D. Pellerin, "Shape-motion based athlete tracking for multi-level action recognition," in *Proc. of AMDO*, 2006, pp. 385–394.

[4] K. Cheung, S. Baker, and T. Kanade, "Shape-from-silhouette across time: Part ii: Applications to human modeling and markerless motion tracking," *Int. J. of Computer Vision*, vol. 63, no. 3, pp. 225–245, 2005.

[5] I. Haritaoglu, D. Harwood, and L. Davis, "W4: Real-time surveillance of people and their activities," *IEEE Trans. on PAMI*, vol. 22, no. 8, pp. 809–830, 2000.
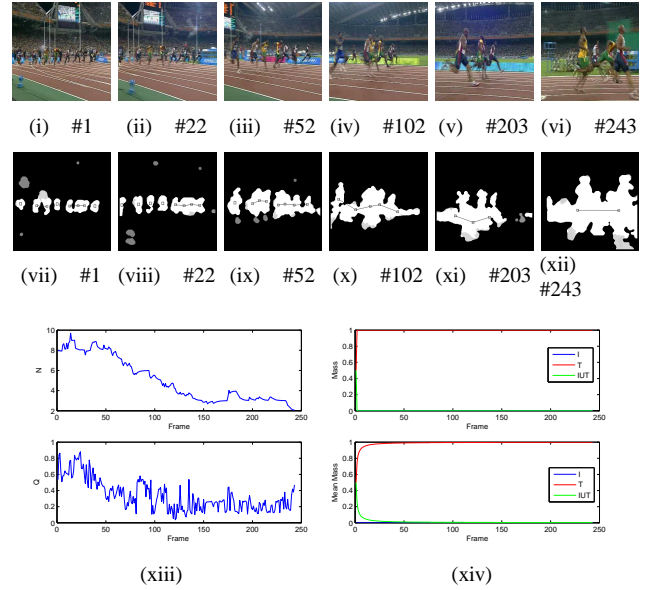
[6] J. Odobez and P. Bouthemy, "Robust multiresolution estimation of parametric motion models," *J. of Vis. Comm. and Image R.*, vol. 6, no. 4, pp. 348–365, 1995.

[7] V. Rabaud and S. Belongie, "Counting crowded moving objects," in *18th IEEE International Conference on Computer Vision and Pattern Recognition*, 2006.

[8] P. J. Figueroa, N. J. Leite, and R. M. L. Barros, "Tracking soccer players aiming their kinematical motion analysis," *CVIU*, vol. 101, no. 2, pp. 122–135, 2006.

[9] P. Smets and R. Kennes, "The Transferable Belief Model," *Art. Intel.*, vol. 66, no. 2, pp. 191–234, 1994.

[10] E. Ramasso, D. Pellerin, C. Panagiotakis, M. Rombaut, G. Tziritas, and W. Lim, "Spatio-temporal information fusion for human action recognition in videos," in *EUSIPCO*, Turkey, 2005.

[11] E. Ramasso, M. Rombaut, and D. Pellerin, "A Temporal Belief Filter improving human action recognition in videos," in *ICASSP*, vol. 2, 2006, pp. 141–144.

[12] D. Dubois and H. Prade, "Representation and combination of uncertainty with belief functions and possibility measures," *Comp. Intel.*, vol. 4, 1988.

[13] P. Smets, "Decision making in the tbm: the necessity of the pignistic transformation." *Int. J. Approx. Reasoning*, vol. 38, no. 2, pp. 133–147, 2005.

[14] R. Kohavi, "The power of decision tables," in *Proc. of European Conference on Machine Learning*, 1995.