Multiple camera people detection and tracking using support integration

Thiago T. Santos, Carlos H. Morimoto

Institute of Mathematics and Statistics University of São Paulo Rua do Matão, 1010 - 05508-090 São Paulo, Brazil

Abstract

This paper proposes a method to locate and track people by combining evidence from multiple cameras using the homography constraint. The proposed method use foreground pixels from simple background subtraction to compute evidence of the location of people on a reference ground plane. The algorithm computes the amount of support that basically corresponds to the "foreground mass" above each pixel. Therefore, pixels that correspond to ground points have more support. The support is normalized to compensate for perspective effects and accumulated on the reference plane for all camera views. The detection of people on the reference plane becomes a search for regions of local maxima in the accumulator. Many false positives are filtered by checking the visibility consistency of the detected candidates against all camera views. The remaining candidates are tracked using Kalman filters and appearance models. Experimental results using challenging data from PETS'06 show good performance of the method in the presence of severe occlusion. Ground truth data also confirms the robustness of the method. *Key words:* people tracking, multiple view integration, video surveillance and monitoring, homography constraint

Preprint submitted to Pattern Recognition Letters

May 30, 2009

1 1. Introduction

This paper presents a multiple camera solution to the problem of tracking 2 people in crowds. Multiple camera views can be used to recover 3D structure 3 information and solve occlusion in crowded environments. Recently, several 4 works have suggested a simpler approach that can be used with a network of 5 sparse uncalibrated cameras based on the homography constraint (1; 2; 3; 4;6 5). The homography constraint establishes that multiple projections of the 7 principal axis of an elongated object using a homography from each camera 8 view q to the ground or reference plane Π intersect at the position of the g object in the reference plane ("ground position" of the object). 10

Kim and Davis (4) use the homography constraint within a particle filter-11 ing framework for people tracking. First, a set of particles that correspond 12 to ground positions is draw from the filter dynamics. Each particle is associ-13 ated with an appearance model (6) to perform people segmentation in each 14 camera view. Once foreground pixels are segmented and classified into sin-15 gle objects (persons), the principal-axis of each person is computed and the 16 homography constraint is used to compute their locations. The main draw-17 back of the system is its requirement that individuals must initially appear 18 as isolated foreground blobs to proper modeling. 19

To detect multiple people using multiple camera views Hu *et al.* (3) use the homography constraint for pairs of cameras. By projecting the principal axis of a person from camera view q to p, the likelihood between two axes from these different views is computed comparing their intersection with a predicted ground position. To compute this point the authors combine single view foreground segmentation with Kalman Filter based tracking. The
likelihood is used to drive the axis correspondence process. The system relies
on individual segmentation, so inter-object occlusion can degrade the axis
location performance.

Eshel and Moses (1) use the homography constraint in several planes par-29 allel to the ground plane, searching for heads in the higher planes. All camera 30 views are mapped using homographies to a reference plane and intensity cor-31 relation is used to detected candidate heads. A nearest neighbor approach 32 is applied to find correspondences along time, producing tracks. In a further 33 step, tracks are combined in individual trajectories by the use of six different 34 measurements to evaluate track overlap, distance, and direction. According 35 to authors, people dressing in similar colors are a main source of false pos-36 itives, a natural drawback from the correlation approach. The cameras are 37 placed at high elevations and the authors report that the performance of the 38 system deteriorates considerably when less than five cameras are used. 39

Fleuret *et al.* (2) use a probabilistic framework to perform simultaneous 40 detection and tracking. Their model is a combination of a simple motion 41 model with an appearance model. The appearance model is composed of an 42 RGB color density and a ground plane occupancy map. In the occupancy 43 map, the ground plane is partitioned into a regular grid and the probability 44 of occupancy of each grid cell is estimated using results from background 45 subtraction. This occupancy model is a conditional distribution between 46 the foreground and the occupied cells configuration. The Viterbi algorithm 47 is used to find the most likely trajectory for each individual and a greedy 48 heuristic is applied to optimize one trajectory after other. For reliable detec-49

tion and location, each person must be seen as an individual blob in at least
one view.

Previous methods for single person segmentation are affected by two main problems. First, partial and total occlusion are common in crowed scenes such as the one in Figure 2b. In places such as airport halls or train stations, people frequently walk in small groups most of the time, causing occlusion in all camera views. Second, when color models are used for segmentation, people dressed with similar colors become another source of problems (3).

The main contribution of this paper is the definition of a novel algo-58 rithm based on the homography constraint that does not rely on single view 59 segmentation of the subjects or previous tracking information. Instead of 60 a segment-then-locate approach, we propose a locate-then-segment approach, 61 integrating available information of all cameras before any detection decision. 62 This paper extends our previous work presented in (5) in several ways. First 63 the people detection method was made more robust to false positives with 64 the introduction of a new filtering algorithm. This paper also introduces a 65 multiple person tracking algorithm based on Kalman filters and appearance 66 models, and more extensive experimental results are presented using ground 67 truth tracking data. 68

Because the system does not require previous object segmentation for people detection, our work has some similarities with the very recent work of Khan and Shah (7). Their work use the homography constraint to fuse foreground likelihood information from multiple views to resolve occlusions and localize people on a reference scene plane. Similar to Eshel and Moses (1), Khan and Shah (7) also rely on multiple planes parallel to the ground to im⁷⁵ prove the robustness of the method. Detection and tracking are performed
⁷⁶ simultaneously by graph cuts segmentation of tracks in the space-time occu⁷⁷ pancy likelihood data.

In our method, multiple view perspective geometry and the homography 78 constraint are applied to collect evidence of people presence from each camera 79 view. Our method elegantly integrates the information of all parallel planes 80 by projecting the foreground directly on the reference plane and accumu-81 lating the evidence from multiple cameras. Occlusion and people detection 82 are solved simultaneously and instantly at each time using the accumulated 83 evidence from all cameras. We have tested the method using very challeng-84 ing data from PETS'06 with good results. The next section describes the 85 method in detail. Experimental results are presented in Section 3. Section 4 86 concludes the paper. 87

88 2. Multiple person detection and tracking

Figure 1 shows a block diagram of our proposed multiple person detection 89 and tracking system. Each static camera q feeds a background subtraction 90 module. The background color distribution for each pixel is modeled us-91 ing mixture of Gaussians. The segmented foreground is used to compute 92 evidence of people presence for each pixel on the reference image Π (floor 93 plane). Our algorithm computes the amount of support that basically cor-94 responds to the "foreground mass" above each pixel. Therefore, pixels that 95 correspond to ground points have more support. Perspective is carefully 96 considered to accurately detect objects near and far away from the cameras. 97 The support computed from each camera view is transformed to the ground 98



Figure 1: Block diagram of the multiple person detection and tracking system.

plane using the appropriate homography. The ground plane accumulates the 99 evidence from all views. People detection is performed by locating regions of 100 local maxima in the ground plane accumulator. Once people candidates are 101 detected, appearance models are computed for each candidate. We have de-102 veloped an efficient algorithm to match the detected candidates with tracked 103 objects. Each tracked object is represented by its appearance model and 104 an associated Kalman filter. Trackers that are assigned to candidates dur-105 ing the matching process are updated. Observations that do not match any 106 tracker are potential new targets, and trackers that do not receive a match 107 are considered lost. 108

109 2.1. Background subtraction

The color distribution for each background pixel in time is modeled as a mixture of Gaussian distributions (8). This Gaussians mixture approach is able to deal with multiple modes on the background color distribution probability.

114

A pixel **x** presents color $f(\mathbf{x})$, represented in rgI space (normalized red,

normalized green and light intensity). Normalized color is less sensitive (compared to RGB space) to small changes in illumination caused by shadows (9).

The color distribution of a pixel is modeled by K Gaussians. The k-117 th Gaussian presents mean vector $\mu_k = \langle \mu_k^r, \mu_k^g, \mu_k^I \rangle$, a diagonal covariance 118 matrix Σ_k and a weight w_k , that correspond to the probability that the pixel 119 has a subclass k. An expectation-maximization (EM) algorithm combined to 120 an agglomerative clustering strategy (10) is applied to estimate K and the 121 mixture parameters of each color distribution. Because the training set is 122 not free from moving objects, the *background distribution* is represented by 123 the Gaussians whose weight w_k is greater than a threshold T_w . 124

Each pixel \mathbf{x}_i is compared against all subclasses in the background mixture model. The pixel is classified as foreground if

$$|f_c(\mathbf{x}_i) - \mu_k^c| > T_b \cdot \sigma_k^c \tag{1}$$

¹²⁷ for all channels c = r, g, I, where T_b is a decision boundary threshold.

Shadows are a common source of artifacts. We use an additional test, based on Wang and Suter (9) work, to perform shadow removal. Let $f^{I}(\cdot)$ denotes the intensity of a pixel in f. If \mathbf{x}_{i} chromaticity fits the pixel r and gmodels and

$$T_{\text{shadow}} \le \frac{f_I(\mathbf{x}_i)}{\mu_k^I} \le 1.0$$

where T_{shadow} is a threshold, then \mathbf{x}_i will be classified as background. The idea is that a background pixel will present just a fraction of its expected intensity value within shadow regions.

135 2.2. Support computation

Let Π be the ground or reference plane, \mathbf{x}^q be a foreground pixel of camera q corresponding to the projection of the point $\mathbf{X} \in \Pi$, and let the pixel relation \mathbf{y} above \mathbf{x} be true iff the foreground pixel \mathbf{y} lies on the half line defined by the ray $\mathbf{x} + \mathbf{u}\mathbf{p}$, and false otherwise, where $\mathbf{u}\mathbf{p}$ is a unit vector pointing to the up direction.

Just for illustration purposes, consider a single person scenario represented by a line segment L. Let $\mathbf{X}_i \in \Pi$ be the bottom end of L, \mathbf{l}^q the projection of L for camera q, and \mathbf{x}_i^q the projection of $\mathbf{X}_i \in \mathbf{l}^q$. Then all pixels $\mathbf{x}_j^q \in \mathbf{l}^q$ such that $i \neq j$, are above \mathbf{x}_i^q . We define support $S(\mathbf{x}_i^q)$ as the number of foreground pixels above \mathbf{x}_i^q .

Notice that $S(\mathbf{x}_i^q)$ can be computed for any \mathbf{x}_i^q regardless of a true cor-146 respondence between \mathbf{x}_i^q and a ground point in Π because only the above 14 relation is used. The vanishing point in the vertical direction can be used to 148 compute the true \vec{up} direction for every pixel \mathbf{x}^q . For a blob corresponding 149 to the segmentation of a person using the background subtraction algorithm, 150 the support of every pixel \mathbf{x}_i^q within the blob can be computed and back-15 projected onto the ground plane. Regions on the ground plane with large 152 local support values are good candidates for the location of a person. 153

154 2.2.1. Perspective normalization

¹⁵⁵ Due to perspective, simple pixel counting to compute $S(\mathbf{x}_i^q)$ is not accu-¹⁵⁶ rate. Figure 2 (b) shows six vertical bars of different lengths. All of them ¹⁵⁷ correspond to the *same* height *h* of the person standing at \mathbf{x}_r^q but at different ¹⁵⁸ locations \mathbf{x}_i^q . Therefore, in order to use support to compute object locations, ¹⁵⁹ the support values must be normalized to compensate for perspective effects.



Figure 2: (a) Perspective transformation for two cameras p and q with projection centers \mathbf{C}^p and \mathbf{C}^q and vanishing points \mathbf{v}^p and \mathbf{v}^q . (b) Perspective correction and height filtering. The bright areas correspond to segmented foreground. The vertical bars correspond to the height of the person standing at \mathbf{x}^q_r seen at different locations \mathbf{x}^q_i .

Using an object of known height h_r as reference, seen by every camera q at \mathbf{x}_r^q , we pre-compute a normalization factor $\eta(\mathbf{x}_i^q)$, for all \mathbf{x}_i^q , that corresponds to the inverse of the height h_r when the reference object is placed at the ground position corresponding to \mathbf{x}_i^q .

For any camera q, let \mathbf{x}_r^q be the position of the reference object with height h_r . Let $\hat{\mathbf{x}}_r^q$ be the projection of \mathbf{x}_r^q onto a parallel plane h_r units far from Π , as shown in Figure 3. Let d(i, j) denote the distance in pixels between any two points (i, j) and assume that $d(\mathbf{x}_r^q, \hat{\mathbf{x}}_r^q)$ is known (the reference height). Then the height $d(\mathbf{x}_i^q, \hat{\mathbf{x}}_i^q)$ of the object when placed at \mathbf{x}_i^q can be estimated using the cross-ratio invariance property of projective geometry (11).

Criminisi *et al.* (11) applied the cross-ratio to find the relation

$$\frac{h_r}{h_q} = 1 - \frac{d(\hat{\mathbf{x}}_r^q, \mathbf{c}_r^q) \, d(\mathbf{x}_r^q, \mathbf{v}^q)}{d(\mathbf{x}_r^q, \mathbf{c}_r^q) \, d(\hat{\mathbf{x}}_r^q, \mathbf{v}^q)} \tag{2}$$

¹⁷¹ between the reference height h_r and the camera height h_q (the distance from ¹⁷² the camera center to the reference plane Π) when the reference object is



Figure 3: Distances for the computation of the perspective normalization factor for the reference position \mathbf{x}_r^q and an arbitrary position \mathbf{x}_i^q . **1** is the ground plane vanishing line (horizon seen by camera q) and \mathbf{v}^q is the vertical vanishing point.

located at \mathbf{x}_r^q . The points \mathbf{c}_r^q and \mathbf{c}_i^q are the projections of \mathbf{x}_r^q and \mathbf{x}_i^q onto the ground plane vanishing line **l**, as seen in Figure 3.

A similar equation can be computed when the reference object is placed at \mathbf{x}_i^q

$$\frac{h_r}{h_q} = 1 - \frac{d(\hat{\mathbf{x}}_i^q, \mathbf{c}_i^q) \, d(\mathbf{x}_i^q, \mathbf{v}^q)}{d(\mathbf{x}_i^q, \mathbf{c}_i^q) \, d(\hat{\mathbf{x}}_i^q, \mathbf{v}^q)}.$$
(3)

Now consider $\alpha(\mathbf{x}_i^q) = d(\mathbf{x}_i^q, \mathbf{v}^q)$ and $\beta(\mathbf{x}_i^q) = d(\mathbf{x}_i^q, \mathbf{c}_i^q)$. Then terms on $\hat{\mathbf{x}}_i^q$ can be rewritten as

$$d(\hat{\mathbf{x}}_{i}^{q}, \mathbf{v}^{q}) = \alpha(\mathbf{x}_{i}^{q}) - \eta(\mathbf{x}_{i}^{q})$$

$$\tag{4}$$

179

$$d(\hat{\mathbf{x}}_{i}^{q}, \mathbf{c}_{i}^{q}) = \beta(\mathbf{x}_{i}^{q}) - \eta(\mathbf{x}_{i}^{q}).$$
(5)

180 Defining

$$\gamma = \frac{d(\hat{\mathbf{x}}_r^q, \mathbf{c}_r^q) \, d(\mathbf{x}_r^q, \mathbf{v}^q)}{d(\mathbf{x}_r^q, \mathbf{c}_r^q) \, d(\hat{\mathbf{x}}_r^q, \mathbf{v}^q)},\tag{6}$$

and using the equality between (2) and (3), it results that:

$$\eta(\mathbf{x}_i^q) = \frac{\alpha(\mathbf{x}_i^q)\beta(\mathbf{x}_i^q)(1-\gamma)}{\alpha(\mathbf{x}_i^q) - \beta(\mathbf{x}_i^q)\gamma}.$$
(7)

The value of $\eta(\mathbf{x}_i^q)$ is pre-computed for each \mathbf{x}_i^q and used as a perspective normalization factor for the computation of support.

184 2.2.2. Bounded support computation

Because objects occlude each other, blobs segmented using background subtraction might be composed of several objects. Large elongated blobs produce large number of false positives due to false high support values. By limiting object heights within an appropriated range $[h_{\min}, h_{\max}]$, the maximum normalized support value is also bounded and the number of false positive candidates is minimized. Small objects with low support values can also be filtered using h_{\min} .

Thus a candidate object for tracking cannot present support below the 192 minimum height h_{\min} or above a maximum h_{\max} . Figure 2 (b) illustrates 193 the idea. Bright areas mark the foreground segmented from camera q. The 194 vertical bar directions are defined by the ground points \mathbf{x}_i^q and the vanishing 195 point \mathbf{v}^q . The bar lengths in pixels correspond to h_{max} . The support of 196 \mathbf{x}_i^q is the amount of foreground pixels along its corresponding bar. Observe 197 that the point \mathbf{x}_1^q does not present any support and that \mathbf{x}_2^q , $\mathbf{x}_3^q,$ \mathbf{x}_4^q and 198 \mathbf{x}_5^q present similar support values. Observe that the line of 3 people under 199 occlusion would cause unrealistically high support values in a large region. 200

The bounded normalized support $S_q(\mathbf{x}_i^q)$ can be computed efficiently for all pixels of a line defined by \mathbf{x}_i^q and \mathbf{v}^q (i.e., a line orthogonal to the ground plane Π) as follows.

Let $\mathbf{s} = \langle \mathbf{x}_1^q, ..., \mathbf{x}_n^q \rangle$ be the line segment obtained by constraining the line by the image frame, as seen in Figure 4. Algorithm 1 computes the support by counting the number of foreground pixels projecting onto \mathbf{x}_i^q and using the perspective normalization factor $\eta(\mathbf{x}_i^q)$ to get the support value in reference units. The maximum support is constrained to filter out objects extending beyond h_{max} .

As an example to better understand the algorithm, consider that at lo-210 cation \mathbf{x}_{280}^q there are 240 foreground pixels above, i.e., F[280] = 240, as seen 211 in Figure 4. According to the pre-computed values of $\eta(\mathbf{x}_{280}^q)$ and h_{max} , the 212 tallest allowed object at location \mathbf{x}_{280}^q would cover up to 120 pixels and reach 213 pixel \mathbf{x}_{160}^q (see line 9 of the algorithm). Since F[160] = 140 (there are 140 214 foreground pixels above \mathbf{x}_{160}^q), there are 100 foreground pixels between \mathbf{x}_{280}^q 215 and \mathbf{x}_{160}^q . This number, normalized by $\eta(\mathbf{x}_{280}^q)$ and bounded, is the support 216 due to the evidence at \mathbf{x}_{280}^q . 21

Background segmentation errors affect the correct computation of an object's support. For example, when people are dressed using colors similar to the background color distribution, parts of their bodies are misdetected. The foreground pixel counting used in Lines 4–8 address this issue and does not constrain support computation to perfect background classification.

Figure 5 shows support results for three different cameras. The figure shows support peaks near people's feet, as expected. Some false foreground detection seen in the top row images are caused by shadows, that produce

Algorithm 1 Algorithm to compute the support $S_q(\mathbf{x}_i^q)$ for all points \mathbf{x}_i^q in

segment s. 1: procedure SUPPORT($\mathbf{s} = \langle \mathbf{x}_1^q, ..., \mathbf{x}_n^q \rangle, h_{\min}, h_{\max}, \eta$) $F[0] \leftarrow 0$ 2: for $i \leftarrow 1, n$ do 3: if \mathbf{x}_i^q is FOREGROUND then 4: $F[i] \leftarrow F[i-1] + 1$ 5:else 6: $F[i] \leftarrow F[i-1]$ 7: end if 8: $j \leftarrow i - h_{\max} \cdot \eta[\mathbf{x}_i^q]$ 9: if j > 0 then 10: $h \leftarrow (F[i] - F[j]) / \eta[\mathbf{x}_i^q]$ 11: else 12: $h \leftarrow F[i]/\eta[\mathbf{x}_i^q]$ 13:end if 14: if $h \ge h_{\min}$ then 15: $S_q(\mathbf{x}_i^q) \leftarrow h$ 16:else 17: $S_q(\mathbf{x}_i^q) \leftarrow 0$ 18:end if 19:end for 20: return S_q 21:22: end procedure



Figure 4: An iteration of Algorithm 1 for (i = 280). Line 9 inspects the pixel \mathbf{x}_{160} , which corresponds to the height of the tallest expected object. Since the value of F[140] = 140, There must be 100 foreground pixels between \mathbf{x}_{160} and \mathbf{x}_{280} .

high support values in regions of the ground plane. Although shadow artifacts
can become an issue in single view processing, multiple view integration is
able to minimize this problem.

229 2.3. Integration of multiple camera views

In the absence of occlusions, the support information computed from a single camera provides sufficient evidence to locate people on the ground plane, though a certain number of false detections and misses might occur. The detection algorithm can be made a lot more robust by combining the evidence from all cameras that see a particular ground region.

For example, in Figure 2 (b), a false ground point \mathbf{x}_3^q has high support but it is unlikely that the same occurs in another camera. In fact, a pair of occluding objects seen in camera q might show as occluding objects for a different camera p iff the objects are along the baseline of the two cameras.

The homography matrix H_q maps ground points \mathbf{x}_i^q in image plane q to ground points \mathbf{X}_i of the ground plane Π according to:

$$\mathbf{x}_i = \mathbf{H}_q \mathbf{x}_i^q. \tag{8}$$

²⁴¹ Using a set of points on the image plane and a set of corresponding points ²⁴² in Π , \mathbb{H}_q can be estimated by a direct linear transformation algorithm (12). ²⁴³ Let $S_q(\mathbf{x}_i^q)$ be the support computed at point \mathbf{x}_i^q for camera q. All support ²⁴⁴ data from Q cameras can be integrated on Π by

$$A(\mathbf{x}_i) = \sum_{q=1}^Q S_q(\mathbf{H}_q^{-1}\mathbf{x}_i).$$
(9)

where A is the *accumulator* image (Figure 6). Objects can be located by segmenting regions of A that present large support values.



Figure 5: (a) input images for the support algorithm. (b) Observe that the support peaks at the ground positions of each person.

A threshold T_S is used to select points $\mathbf{X}_i \in \Pi$ presenting good support values. The threshold parameter at $\mathbf{X}_i \in \Pi$ takes into consideration h_{\min} and the number of cameras able to see that location. Points of local maxima are computed by a mean-shift procedure. Mean-shift blurring process (13) moves data points in the gradient direction of a smoothed version of the original function. Applied to A, the process integrates the support information within a neighborhood of \mathbf{X}_i .

254

Let G be the set of found local maxima points. Points $\mathbf{X}_i \in G$ cor-

respond to real people locations and some false-positives. Main sources of false-positives are severe occlusion in all views and people aligned in the baseline of a pair of cameras. The idea to filter the false-positives is to select a subset of G that, under total occlusion relations, is able to "explain" the occurrence of the remaining points.

Points in G are labeled UNSELECTED and inserted in a priority queue 260 ordered by $A(\mathbf{X}_i)$. We pop the queue, marking the current point \mathbf{X}_i as SE-26 LECTED. Then we visit all the points \mathbf{X}_j that are occluded by \mathbf{X}_i . If \mathbf{X}_j is 262 UNSELECTED and it is occluded by a SELECTED point in all views, it will 263 be labeled COVERED and removed from the queue. We repeat this proce-264 dure until no more UNSELECTED points are available. SELECTED points are 265 returned as people location candidates and will be further used as measure-266 ments by the tracking module. This procedure ensures that the removed 26 false-positives are fully justified as spurious interactions from evidences of 268 people in other locations. 269

270 2.4. Object Tracking

Our system tracks multiple objects simultaneously using one Kalman Filter per object. A tracked object (person) is represented by a multi-view appearance model. The model consists of two RGB color histograms for each camera view, corresponding to the top an bottom parts of the object (shirt and pants). Each model also keeps a foreground and occlusion mask for each camera. The color histograms, foreground, and occlusion masks are updated at every frame.

Before updating the tracker at every new frame t, appearance models for the detected target candidates (called the observation appearance mod-



Figure 6: Multi-view integration for 3 cameras. Homographies are used to warp support from the original camera view to the floor plane Π . The accumulated support $A(\mathbf{x}_i)$ peaks on true object positions.

els) are build using the list of candidate positions computed as described
in previously. A bounding box for each camera view is computed from the
position and estimated height (support) of the candidate object. The RGB
color histograms, foreground, and occlusion masks are computed using such
bounding boxes.

To efficiently determine the assignment of observations to targets all possible assignments we have developed the following greedy algorithm.

First candidate positions z_i are paired with all trackers T_j that expects the tracked object to be at a vicinity of $_z_i$. All such pairs are inserted in a priority queue according to the probability $p(z_i|x_j, \sigma_j)$, where z_i is the observation position on the ground plane and x_j and σ_j are respectively the state and covariance matrix of the Kalman Filter T_j .

Next the first pair of the queue is popped and their appearance models 292 are used to test if the observation actually matches the tracked object. An 293 observation matches an object iff there is good similarity between their color 294 models. Color similarity is computed using histogram intersection. In case 295 the tracker is updated using the matched observation, the object appearance 296 model is also updated using the observation appearance model and a learn-29 ing factor alpha as follows. Let $H_{q,t}[b]$ be the histogram value for bin b in 298 the a color model of camera q at frame t and let H^o be the corresponding 299 observation model. Then 300

$$H_{q,t+1}[b] = (1 - \alpha)H_{q,t}[b] + \alpha H_{q,t+1}^{o}[b]$$
(10)

Observation z_i that are matched are marked as USED, so no other tracker 301 will be updated using z_i . The process continues until the queue is empty. 302 The greedy algorithm might not assign all observations to all trackers. Ob-303 servations that are not assigned to a tracker correspond to potential new 304 objects so a new tracker is created. Each tracker T_j keeps a counter to reg-305 ister the number of successful assignments, and a flag. Upon creation new 306 trackers receive a NEW flag and their appearance models initialized to the 30 observation appearance models. 308

After the counter registers a large enough number of assignments, the tracker flag is updated to ON. At this moment, the tracker is assumed to be following a real subject. If a tracker is not assigned to any observation, its flags is updated to LOST. A LOST tracker is updated using the Kalman prediction and its covariance matrix is increased to enhance the chances of the tracker to find a match in the next frame. A tracker that keeps a LOST flag for a long time is finished and removed from the list of trackers. Trackers presenting the ON flag have priority on the assignment queue and
LOST trackers have priority over NEW ones.

318 3. Results

The system was tested using the S7 dataset from the PETS 2006 Bench-319 mark Data (14). This dataset presents video recorded at Victoria Station in 320 London, UK. Video from three cameras was used, demonstrating that just 32 a few cameras are enough to produce good detection and tracking results. 322 We used half of S7 frame sequence in our tests (the last 1500 frames of the 323 original 3000 sequence - about 1 minute of video). The sequence presents 324 22 individuals walking in a hall. About 1/3 of the hall area is covered by 325 three cameras. The baseline of the two cameras that cover the remaining 326 area crosses the entire hall, creating severe occlusion situations. 32

Image points were manually selected to compute the vanishing points of 328 each camera and the appropriate homography matrix to the ground plane 329 Π . The height of a person was used to define the reference height unit. The 330 allowed height range was set to [0.6, 1.1] units (that is 60% to 110% of the 33 reference man's height). An unit flat kernel of width 19 pixels was applied 332 in the mean-shift local maxima detection procedure (1 pixel ~ 2 cm in the 333 reference ground plane image). Trajectories from the tracking module shorter 334 than 50 frames (about 2 s) are considered false-positives and removed. 335

336 3.1. Object Detection

Figures 7 and 8 show results for two situations presenting occlusion cases. The first row displays the floor plane square texture pattern and the detected object positions. These points are classified as people's ground points and





Figure 7: Local maxima corresponds to location of people on the reference ground plane (marked with dots). The homographies H_q are used to map the people's ground points back to each camera view.

are shown as red dots in the next row. Homographies are used to map the
ground points back to each camera view.

The subjects of interest are the people visible on the floor plane diagram in the first row of Figure 7. Frame 3300 in Figure 8 shows an example of occlusion under three views. The proposed system is able to detect each individual successfully.

346 3.2. Tracking

Ground-truth was manually created to evaluate tracking results. The position of each individual was manually annotated for 150 frames, 10 frames apart for the 1500 frames of the S7 PETS sequence. Consistent labeling was



Figure 8: Another example from PETS'06 dataset. Frame 3300 presents occlusion in all camera views but the system could accurately find the right people location.

	PETS 2006 S07
Number of Trajectories	22
Found tracks	30
Trajectory Recall	100.00%
Trajectory Precision	96.67%
Tracks per Trajectory	1.3182

Table 1: Tracks found by the tracking procedure compared to ground truth people trajectories.

associated to each person. Table 1 summarizes the results. All 22 subjects 350 were successfully associated to one or more tracks produced by the system. 351 Only one of the tracks does not match any subject. Ideally, one tracker 352 should be associated to one person for the whole sequence. The proposed 353 system produced an average of 1.32 tracks per trajectory, which corresponds 354 to few errors during tracking. There was only one track exchange amongst all 355 trackers for the whole sequence that took place between two near individuals, 356 seen only by 2 cameras, in occlusion and aligned to the cameras baseline. 357

Figure 9 shows the root mean square deviation between the estimated trajectories and the ground truth positions for each subject. The largest deviation was about 50 cm and its associated to a running man in the video sequence (subject 14). Figure 10 displays the estimated and ground truth trajectory for subject 19. This subject crosses the entire hall and is occluded by other people several times.



Figure 9: Root mean square deviation for PETS 2006 S07 sequence.



Figure 10: Trajectory for subject 19. The subject was occluded several times along the trajectory. There are foreground misdetection at some points, caused by color similarity between his clothes and the background. The baseline between cameras 1 and 3 is marked as a dashed gray line.

4. Conclusions

novel method to locate people on the ground plane using multiple cam-365 era views was presented. The main advantage of the method is that it does 366 not require initial people segmentation or tracking. The robustness of the 36 method is due to the accumulation of support from all cameras. The support 368 of a candidate object location is defined as the amount of foreground pix-360 els above that location. Therefore, pixels that correspond to ground points 370 have more support. The support is normalized to compensate for perspec-371 tive effects and accumulated on the reference plane for all camera views. The 372 detection of people on the reference plane becomes a search for regions of 373 local maxima in the accumulator. The paper also introduces a filtering algo-374 rithm that eliminates many false positives by checking the consistency of the 375 location against the remaining objects for all camera views. The remaining 376 candidates are tracked using Kalman filters and appearance models. Chal-37 lenging sequences from PETS'2006 were used to test the system and show its 378 robustness to severe occlusion situations using just 3 sparse cameras. Ground 370 truth data also confirms the tracking accuracy of the method. 380

Future work includes further experimentation in other crowded scenarios and trajectory analysis for event detection.

383 Acknowledgments

T. T. Santos acknowledges support from Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES – grant BEX 2686/06). T. T. Santos and C. H. Morimoto acknowledge financial support from Fundação de Amparo à Pesquisa do Estado de São Paulo (FAPESP).

388 References

- [1] R. Eshel, Y. Moses, Homography based multiple camera detection and tracking of people in a dense crowd, in: Proceedings of 2008 IEEE Conference on Computer Vision and Pattern Recognition (CVPR 2008), Los
 Alamitos, CA, USA, 2008, pp. 1–8. doi:10.1109/CVPR.2008.4587539.
- F. Fleuret, J. Berclaz, R. Lengagne, P. Fua, Multicamera people
 tracking with a probabilistic occupancy map, Pattern Analysis and
 Machine Intelligence, IEEE Transactions on 30 (2) (2008) 267–282.
 doi:10.1109/TPAMI.2007.1174.
- [3] W. Hu, M. Hu, X. Zhou, T. Tan, J. Lou, S. Maybank, Principal axisbased correspondence between multiple cameras for people tracking,
 Pattern Analysis and Machine Intelligence, IEEE Transactions on 28 (4)
 (2006) 663-671. doi:10.1109/TPAMI.2006.80.
- [4] K. Kim, L. Davis, Multi-camera tracking and segmentation of occluded
 people on ground plane using search-guided particle filtering, in: Proceedings of 9th European Conference on Computer Vision (ECCV'06),
 Vol. 3953, Graz, Austria, 2006, pp. 98–109. doi:10.1007/11744078_8.
- [5] T. T. Santos, C. H. Morimoto, People detection under occlusion in multiple camera views, in: Proceedings of XXI Brazilian Symposium on Computer Graphics and Image Processing (SIBGRAPI '08), IEEE Computer
 Society, Los Alamitos, 2008, pp. 53–60. doi:10.1109/SIBGRAPI.2008.25.
- ⁴⁰⁹ [6] A. Senior, A. Hampapur, Y.-L. Tian, L. Browna, S. Pankantia, R. Bolle,

Appearance models for occlusion handling, Image and Vision Computing 24 (11) (2006) 1233–1243.

- [7] S. Khan, M. Shah, Tracking multiple occluding people by localizing on
 multiple scene planes, Pattern Analysis and Machine Intelligence, IEEE
 Transactions on 31 (3) (2009) 505–519. doi:10.1109/TPAMI.2008.102.
- [8] C. Stauffer, W. Grimson, Adaptive background mixture models for realtime tracking, in: Proceedings of 1999 IEEE Conference on Computer
 Vision and Pattern Recognition (CVPR'99), Vol. 2, Los Alamitos, CA,
 USA, 1999, pp. 246–252.
- [9] H. Wang, D. Suter, A re-evaluation of mixture of gaussian background
 modeling, in: Proceedings of 30th IEEE International Conference on
 Acoustics, Speech, and Signal Processing (ICASSP 2005), Vol. 2, 2005,
 pp. 1017–1020.
- [10] C. A. Bouman, Cluster: An unsupervised algorithm for modeling Gaussian mixtures, available from http://www.ece.purdue.edu/~bouman
 (April 1997).
- [11] A. Criminisi, I. D. Reid, A. Zisserman, Single view metrology, International Journal of Computer Vision 40 (2) (2000) 123–148.
- [12] R. Hartley, A. Zisserman, Multiple View Geometry in Computer Vision,
 Cambridge University Press, 2004.
- [13] Y. Cheng, Mean shift, mode seeking, and clustering, Pattern Analysis
 and Machine Intelligence, IEEE Transactions on 17 (8) (1995) 790–799.
 doi:10.1109/34.400568.

- 433 [14] D. Thirde, L. Li, J. Ferryman, Overview of the PETS2006 challenge,
- in: Proceedings of 9th IEEE International Workshop on Performance
- ⁴³⁵ Evaluation of Tracking and Surveillance (PETS 2006), New York, USA,
- 436 2006, pp. 47–50.