

# Feature Compression: A Framework for Multi-View Multi-Person Tracking in Visual Sensor Networks

Serhan Coşar\*, Müjdat Çetin

*Sabancı University, Faculty of Engineering and Natural Sciences, Orta Mahalle, Üniversite  
Caddesi No: 27 34956 Tuzla-İstanbul, Turkey*

---

## Abstract

Visual sensor networks (VSNs) consist of image sensors, embedded processors and wireless transceivers which are powered by batteries. Since the energy and bandwidth resources are limited, setting up a tracking system in VSNs is a challenging problem. In this paper, we present a framework for human tracking in VSNs. The traditional approach of sending compressed images to a central node has certain disadvantages such as decreasing the performance of further processing (i.e., tracking) because of low quality images. Instead, in our method, each camera performs feature extraction and obtains likelihood functions. By transforming to an appropriate domain and taking only the significant coefficients, these likelihood functions are compressed and this new representation is sent to the fusion node. An appropriate domain is selected by performing a comparison between well-known transforms. We have applied our method for indoor people tracking and demonstrated the superiority of our system over the traditional approach.

*Keywords:* Visual sensor networks, Camera networks, Human tracking, Communication constraints, Compressing likelihood functions

---

---

\*Corresponding author. Tel: +90 216 4830000-2117, Fax: +90 216 483-9550.  
*Email addresses:* [serhancosar@sabanciuniv.edu](mailto:serhancosar@sabanciuniv.edu) (Serhan Coşar),  
[mcetin@sabanciuniv.edu](mailto:mcetin@sabanciuniv.edu) (Müjdat Çetin)

## 1. Introduction

With the birth of wireless sensor networks, new applications are enabled by large-scale networks of small devices capable of (i) measuring information from the physical environment, such as temperature, pressure, etc., (ii) performing simple processing on the extracted data, and (iii) transmitting the processed data to remote locations by also considering the limited resources such as energy and bandwidth. More recently, the availability of inexpensive hardware such as CMOS cameras that are able to capture visual data from the environment has supported the development of Visual Sensor Networks (VSNs), i.e., networks of wirelessly interconnected devices that acquire video data.

Using a camera in a wireless network leads to unique and challenging problems that are more complex than the traditional wireless sensor networks might have. For instance, most sensors provide measurements of temporal signals that represent physical quantities such as temperature. On the other hand, at each time instant image sensors provide a 2D set of data points, which we see as an image. This richer information content increases the complexity of data processing and analysis. Performing complex tasks, such as tracking, recognition, etc., in a communication-constrained VSN environment is extremely challenging. With a data compression perspective, the common approach is to compress images and collect them in a central unit to perform the tasks of interest. In this strategy, the main goal is to focus on low-level communication. The communication load is decreased by compressing the raw data without regard to the final inference goal based on the information content of the data. Since such a strategy will affect the quality of the transmitted data, it may decrease the performance of further inference tasks. In this paper, we propose a different strategy for decreasing the communication that is better matched to problems with a defined final inference goal, which, in the context of this paper, is tracking.

There has been some work proposed for solving the problems mentioned above.

31 To minimize the amount of data to be communicated, in some methods simple  
32 features are used for communication. For instance, 2D trajectories are used  
33 in [1]. In [2], 3D trajectories together with color histograms are used. Hue  
34 histograms along with 2D position are used in [3]. Moreover, there are decen-  
35 tralized approaches in which cameras are grouped into clusters and tracking is  
36 performed by local cluster fusion nodes. This kind of approaches have been  
37 applied to the multi-camera target tracking problem in various ways [4, 5, 6].  
38 For a nonoverlapping camera setup, tracking is performed by maximizing the  
39 similarity between the observed features from each camera and minimizing the  
40 long-term variation in appearance using graph matching at the fusion node [4].  
41 For an overlapping camera setup, a cluster-based Kalman filter in a network  
42 of wireless cameras is proposed in [5, 6]. Local measurements of the target ac-  
43 quired by members of the cluster are sent to the fusion node. Then, the fusion  
44 node estimates the target position via an extended Kalman filter, relating the  
45 measurements acquired by the cameras to the actual position of the target by  
46 nonlinear transformations.

47

48 Previous works proposed for VSNs have some handicaps. The methods in  
49 [1, 2, 3] that use simpler features may be capable of decreasing the commu-  
50 nication, but they are not capable of maintaining robustness. For the sake  
51 of bandwidth constraints, these methods choose to change the features from  
52 complex and robust to simpler but not so effective ones. As in the methods  
53 proposed in [4, 5, 6], performing local processing and collecting features to the  
54 fusion node may not satisfy the bandwidth requirements in a communication-  
55 constrained VSN environment. In particular, depending on the size of image  
56 features and the number of cameras in the network, even collecting features to  
57 the fusion node may become expensive for the network. In such cases, further  
58 approximations on features are necessary. An efficient approach that reduces  
59 the bandwidth requirements without significantly decreasing the quality of im-  
60 age features is needed.

61

62 In this paper, we propose a framework that is suitable for energy and band-  
63 width constraints in VSNs. It is capable of performing multi-person tracking  
64 without significant performance loss. Our method is a decentralized tracking  
65 approach in which each camera node in the network performs feature extraction  
66 by itself and obtains image features (likelihood functions). Instead of directly  
67 sending likelihood functions to the fusion node, a block-based compression is  
68 performed on likelihoods by transforming each block to an appropriate domain.  
69 Then, in this new representation we only take the significant coefficients and  
70 send them to the fusion node. Hence, multi-view tracking can be performed  
71 without overloading the network. The main contribution of this work is the  
72 idea of performing goal-directed compression in a VSN. In the tracking context,  
73 this is achieved by performing local processing at the nodes and compressing  
74 the resulting likelihood functions which are related to the tracking goal, rather  
75 than compressing raw images. To the best of our knowledge, compression of  
76 likelihood functions computed in the context of tracking in a VSN has not been  
77 proposed in previous work.

78

79 We have used our method within the context of a well-known multi-camera  
80 human tracking algorithm [7]. We have modified the method in [7] to obtain  
81 a decentralized tracking algorithm. In order to choose an appropriate domain  
82 for likelihood functions, we have performed a comparison between well-known  
83 transforms. A traditional approach in camera networks is transmitting com-  
84 pressed images. Both by qualitative and quantitative results, we have shown  
85 that our method is better than the traditional approach of sending compressed  
86 images and can work under VSN constraints without degrading the tracking  
87 performance.

88

89 In Section 2, how we integrate multi-view information in our decentralized ap-  
90 proach is described. Section 3 presents our feature compression framework in  
91 detail and contains a comparison of various domains for likelihood representa-  
92 tion. Experimental setup and results are given in Section 4. Finally in Section

93 5, we conclude and suggest a number of directions for potential future work.

## 94 **2. Multi-Camera Integration**

### 95 *2.1. Decentralized Tracking*

96 In a traditional setup of camera networks, which we call centralized tracking,  
97 each camera acquires an image and sends this raw data to a central unit. In  
98 the central unit, multi-view data are collected, relevant features are extracted  
99 and combined, finally, using these features, the positions of the humans are  
100 estimated. Hence, integration of multi-view information is done in raw-data  
101 level by pooling all images in a central unit. The presence of a single global  
102 fusion center leads to high data-transfer rates and the need for a computation-  
103 ally powerful machine, thereby, to a lack of scalability and energy efficiency.  
104 Compressing raw image data may decrease the communication in the network,  
105 but since the quality of images drops, it might also decrease the tracking per-  
106 formance. For this reason, centralized trackers are not very appropriate for use  
107 in VSN environments. In decentralized tracking, there is no central unit that  
108 collects all raw data from the cameras. Cameras are grouped into clusters and  
109 nodes communicate with their local cluster fusion nodes only [8]. Communi-  
110 cation overhead is reduced by limiting the cooperation within each cluster and  
111 among fusion nodes. After acquiring the images, each camera extracts useful  
112 features from the images it has observed and sends these features to the local  
113 fusion node. Using the multi-view image features, tracking is performed in the  
114 local fusion node. Hence, we can say that in decentralized tracking, multi-view  
115 information is integrated in feature-level by combining the features in small clus-  
116 ters. The decentralized approaches fits very well to VSNs in many aspects. The  
117 processing capability of each camera is utilized by performing feature extraction  
118 at camera-level. Since cameras are grouped into clusters, the communication  
119 overhead is reduced by limiting the cooperation within each cluster and among  
120 fusion nodes. In other words, by a decentralized approach, feature extraction  
121 and communication are distributed among cameras in clusters, therefore, effi-

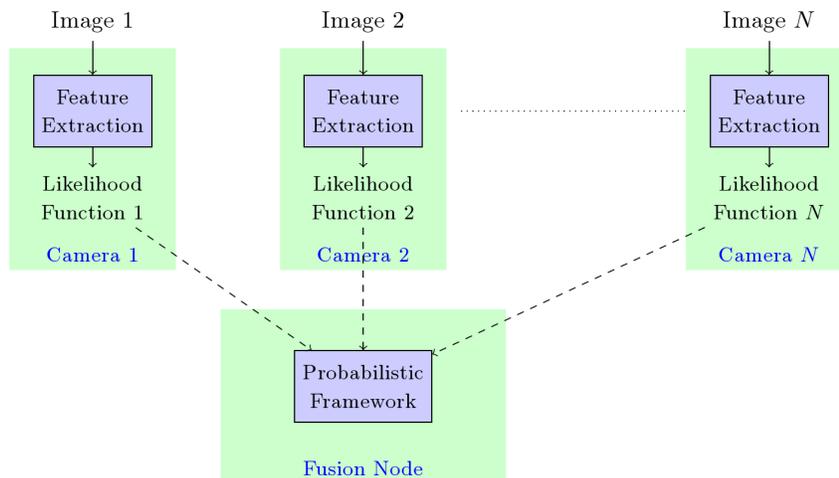


Figure 1: The flow diagram of a decentralized tracker using a probabilistic framework.

122 cient estimation can be performed.

123

124 Modeling the dynamics of humans in a probabilistic framework is a common  
 125 perspective of many multi-camera human tracking methods [7, 9, 10, 11]. In  
 126 tracking methods based on a probabilistic framework, data and/or extracted fea-  
 127 tures are represented by likelihood functions,  $p(y|x)$  where  $y \in R^d$  and  $x \in R^m$   
 128 are the observation and state vectors, respectively. In other words, for each  
 129 camera, a likelihood function is defined in terms of the observations obtained  
 130 from its field of view. In centralized tracking, of course, the likelihood functions  
 131 are computed after collecting the image data of each camera at the central unit.  
 132 For a decentralized approach, since each camera node extracts local features  
 133 from its field of view, these likelihood functions can be evaluated at the camera  
 134 nodes and they can be sent to the fusion node. Then, in the fusion node the  
 135 likelihoods can be combined and tracking can be performed in the probabilistic  
 136 framework. A flow diagram of the decentralized approach is illustrated in Fig-  
 137 ure 1. Following this line of thought, we have converted the tracking approach  
 138 described in Section 2.2 to a decentralized tracker as explained in Section 2.3.

139 *2.2. Multi-Camera Tracking Algorithm*

140 In this section we describe the tracking method of [7], as we apply our pro-  
 141 posed approach within in the context of this method in this paper. In [7],  
 142 the visible part of the ground plane is discretized into a finite number  $G$  of  
 143 regularly spaced 2D locations. Let  $\mathbf{L}_t = (L_t^1, \dots, L_t^{N^*})$  be the locations of in-  
 144 dividuals at time  $t$ , where  $N^*$  stands for the maximum allowable number of  
 145 individuals. Given  $T$  temporal frames from  $C$  cameras,  $\mathbf{I} = (\mathbf{I}_1, \dots, \mathbf{I}_T)$  where  
 146  $\mathbf{I}_t = (I_t^1, \dots, I_t^C)$ , the goal is to maximize the posterior conditional probability:

$$P(\mathbf{L}^1 = \mathbf{l}^1, \dots, \mathbf{L}^{N^*} = \mathbf{l}^{N^*} | \mathbf{I}) = P(\mathbf{L}^1 = \mathbf{l}^1 | \mathbf{I}) \prod_{n=2}^{N^*} P(\mathbf{L}^n = \mathbf{l}^n | \mathbf{I}, \mathbf{L}^1 = \mathbf{l}^1, \dots, \mathbf{L}^{n-1} = \mathbf{l}^{n-1}) \quad (1)$$

147 where  $\mathbf{L}^n = (L_1^n, \dots, L_T^n)$  is the trajectory of person  $n$ . Simultaneous optimiza-  
 148 tion of all the  $L^i$ 's would be intractable. Instead, one trajectory after the other  
 149 is optimized.  $\mathbf{L}^n$  is estimated by seeking the maximum of the probability of  
 150 both the observations and the trajectory ending up at location  $k$  at time  $t$ :

$$\Phi_t(k) = \max_{l_1^n, \dots, l_{t-1}^n} P(\mathbf{I}_1, L_1^n = l_1^n, \dots, \mathbf{I}_t, L_t^n = k) \quad (2)$$

151 Under a hidden Markov model, the above expression turns into the classical  
 152 recursive expression:

$$\Phi_t(k) = \underbrace{P(\mathbf{I}_t | L_t^n = k)}_{\text{Appearance model}} \max_{\tau} \underbrace{P(L_t^n = k | L_{t-1}^n = \tau)}_{\text{Motion model}} \Phi_{t-1}(\tau) \quad (3)$$

153 The motion model  $P(L_t^n = k | L_{t-1}^n = \tau)$  is a distribution into a disc of limited  
 154 radius and center  $\tau$ , which corresponds to a loose bound on the maximum speed  
 155 of a walking human.

156

157 From the input images  $\mathbf{I}_t$ , by using background subtraction, foreground bi-  
 158 nary masks,  $\mathbf{B}_t$ , are obtained. Let the colors of the pixels inside the blobs are  
 159 denoted as  $\mathbf{T}_t$  and  $X_k^t$  be a Boolean random variable denoting the presence of  
 160 an individual at location  $k$  of the grid at time  $t$ . It is shown in [7] that the

161 appearance model in Eq. 3 can be decomposed as:

$$\begin{aligned}
 \overbrace{P(\mathbf{I}_t | L_t^n = k)}^{\text{Appearance model}} &\propto \underbrace{P(L_t^n = k | X_k^t = 1, \mathbf{T}_t)}_{\text{Color model}} \underbrace{P(X_k^t = 1 | \mathbf{B}_t)}_{\text{Ground plane occupancy}} \quad (4)
 \end{aligned}$$

162

163 In [7], humans are represented as simple rectangles and these rectangles are used  
 164 to create synthetic ideal images that would be observed if people were at given  
 165 locations. Within this model, the ground plane occupancy is approximated by  
 166 measuring the similarity between ideal images and foreground binary masks.

167

168 Let  $T_t^c(k)$  denote the color of the pixels taken at the intersection of the fore-  
 169 ground binary mask,  $B_t^c$ , from camera  $c$  at time  $t$  and the rectangle  $A_k^c$  corre-  
 170 sponding to location  $k$  in that same field of view. Say we have the reference color  
 171 distributions (histograms) of the  $N^*$  individuals present in the scene,  $\mu_1^c, \dots, \mu_{N^*}^c$ .

172 The color model of person  $n$  in Eq. 4 can be expressed as:

$$\begin{aligned}
 \overbrace{P(L_t^n = k | X_k^t = 1, \mathbf{T}_t)}^{\text{Color model}} &\propto P(\mathbf{T}_t | L_t^n = k) = P(T_t^1(k), \dots, T_t^C(k) | L_t^n = k) \\
 &= \prod_{c=1}^C P(T_t^c(k) | L_t^n = k) \quad (5)
 \end{aligned}$$

173 In [7], by assuming the pixels whose colors are represented by  $T_t^c(k)$  are in-  
 174 dependent,  $P(T_t^c(k) | L_t^n = k)$  is evaluated by a product of the marginal color  
 175 distribution  $\mu_n^c$  at each pixel,  $P(T_t^c(k) | L_t^n = k) = \prod_{r \in T_t^c(k)} \mu_n^c(r)$ . In this ap-  
 176 proach, a patch with constant color intensity corresponding to the the mode  
 177 of the color distribution would be most likely. Hence, this approach may  
 178 fail to capture the statistical color variability represented by the full proba-  
 179 bility density function estimated from a spatial patch. Instead, we represent  
 180  $P(T_t^c(k) | L_t^n = k)$  by comparing the observed and reference color distribu-  
 181 tions, which is a well known approach used in many computer vision methods  
 182 [12, 13, 14]. In particular, we compare the estimated color distribution (his-  
 183 togram) of the pixels in  $T_t^c(k)$  and the color distribution  $\mu_n^c$  with a distance  
 184 metric -  $P(T_t^c(k) | L_t^n = k) = \exp(-S(H_t^{c,k}, \mu_n^c))$  where  $H_t^{c,k}$  denotes the his-  
 185 togram of the pixels in  $T_t^c(k)$  and  $S(\cdot)$  is a distance metric. As a distance

186 metric, we use the Bhattacharya coefficient between two distributions. In this  
187 way, we can evaluate the degree of match between the intensity distribution of  
188 an observed patch and the reference color distribution.

189

190 By performing a global search with dynamic programming using Eq. 3, the  
191 trajectory of each person can be estimated.

### 192 2.3. Decentralized Version of the Tracking Algorithm

193 From the above formulation, we can see that there are two different likeli-  
194 hood functions defined in the method. One is the ground plane occupancy map  
195 (GOM),  $P(X_k^t = 1 | \mathbf{B}_t)$ , approximated using the foreground binary masks. The  
196 other is the ground plane color map (GCM),  $P(L_t^n = k | X_k^t = 1, \mathbf{T}_t)$ , which is a  
197 multi-view color likelihood function defined for each person individually. This  
198 map is obtained by combining the individual color maps,  $P(T_t^c(k) | L_t^n = k)$ ,  
199 evaluated using the images each camera acquired. Since foreground binary  
200 masks are simple binary images that can be easily compressed by a lossless  
201 compression method, they can be directly sent to the fusion node without over-  
202 loading the network. Therefore, we keep these binary images as in the original  
203 method and GOM is evaluated at the fusion node. In our framework, we eval-  
204 uate GCM in a decentralized way (as presented in Figure 1): At each camera  
205 node ( $c = 1, \dots, C$ ), the local color likelihood function for the person of interest  
206 ( $P(T_t^c(k) | L_t^n = k)$ ) is evaluated by using the image acquired from that camera.  
207 Then, these likelihood functions are sent to the fusion node. At the fusion node,  
208 these likelihood functions are integrated to obtain the multi-view color likeli-  
209 hood function (GCM) (Eq. 5). By combining GCM and GOM with the motion  
210 model, the trajectory of the person of interest is estimated at the fusion node  
211 using dynamic programming (Eq. 3). The whole process is run for each person  
212 in the scene.

213

214 Fusion node selection and sensor resource management (sensor tasking) is out of  
215 scope of this paper. We have assumed that one of the camera nodes, relatively

216 more powerful one, has been selected as the fusion node.

### 217 **3. Feature Compression Framework**

#### 218 *3.1. Compressing Likelihood Functions*

219 The bandwidth required for sending local likelihood functions depends on  
220 the size of likelihoods (i.e., the number of "pixels" in a 2D likelihood function)  
221 and the number of cameras in the network. To make the communication in the  
222 network feasible, we propose a feature compression framework. In our frame-  
223 work, similar to image compression, we compress the likelihood functions by  
224 transforming them to a proper domain and keeping only the significant coef-  
225 ficients, assuming significant parts of the likelihood functions are sufficient for  
226 performing tracking. At each camera node, we first split the likelihood function  
227 into blocks. Then, we transform each block to a proper domain and take only  
228 the significant coefficients in the new representation. Instead of sending the  
229 function itself, we send this new representation of each block. In this way, we  
230 reduce the communication in the network.

231

232 Mathematically, we have the following linear system:

$$y_c^b = A \cdot x_c^b \quad (6)$$

233 where  $y_c^b$  and  $x_c^b$  represent the  $b$ th block of the likelihood function of camera  $c$   
234 (for a person of interest in a particular time instant,  $P(T_t^c(k)|L_t^n = k)$  in Eq. 5)  
235 and its representation, respectively, and  $A$  is the domain we transform  $y_c^b$  to. In  
236 most of the compression methods, the matrix  $A$  is chosen to be a unitary matrix.  
237 Hence, we can obtain  $x_c^b$  by multiplying  $y_c^b$  with the Hermitian transpose of  $A$ :

$$x_c^b = A^* \cdot y_c^b \quad (7)$$

238 Figure 2 illustrates our likelihood compression scheme.

239

240 Notice that in our feature compression framework, we do not require the use

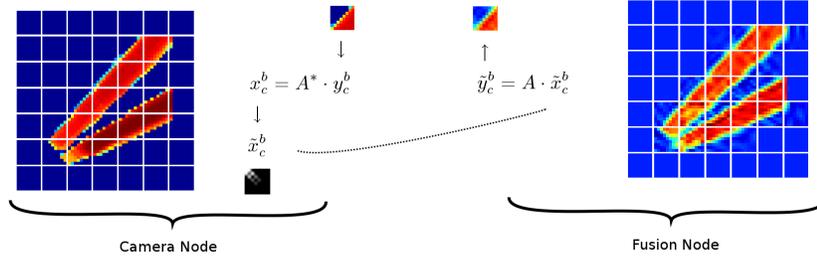


Figure 2: Our Likelihood compression scheme. On the left, there is a local likelihood function ( $P(T_t^c(k)|L_t^n = k)$  in Eq. 5). First, we split the likelihood into blocks, then we transform each block to the domain represented by matrix  $A$  and obtain the representation  $x_c^b$ . We only take significant coefficients in this representation and obtain a new representation  $\tilde{x}_c^b$ . For each block, we send this new representation to fusion node. Finally, by reconstructing each block we obtain the whole likelihood function on the right.

241 of specific image features or likelihood functions. The only requirement is that  
 242 the tracking method should be based on a probabilistic framework, which is a  
 243 common approach for modeling the dynamics of humans. Hence, our frame-  
 244 work is a generic framework that can be used with many probabilistic tracking  
 245 algorithms in a VSN environment.

246

247 In all camera nodes and fusion nodes, the matrix  $A$  is common, therefore, at the  
 248 fusion node, likelihood functions of each camera can be reconstructed simply by  
 249 multiplying the new representation with the matrix  $A$ . In general, this may  
 250 require an offline coordination step to decide the domain that is matched with  
 251 the task of interest. In the next subsection, we go through the question of which  
 252 domain should be selected in Eq. (6).

### 253 3.2. A Proper Domain for Compression

254 By sending the compressed likelihoods to the fusion node, our goal is to  
 255 decrease the communication in the network without affecting the tracking per-  
 256 formance significantly. On one hand, we want to send less coefficients, on the  
 257 other hand, we do not want to decrease the quality of the likelihoods, i.e., we  
 258 want to have small reconstruction error. For this reason, we need to select a

259 domain that is well-matched to the likelihood functions, providing the oppor-  
260 tunity to accurately reconstruct the likelihoods back using a small number of  
261 coefficients.

262

263 Image compression using transforms is a mature research area. Numerous trans-  
264 forms such as the discrete cosine transform (DCT), the Haar transform, symm-  
265 lets, coiflets have been proposed and proven to be successful [15, 16, 17]. DCT  
266 is a well-known transform that has the ability to analyze non-periodic signals.  
267 Haar wavelet is the first known wavelet basis that consists of orthonormal func-  
268 tions. In wavelet theory, *number of vanishing moments* and *size of support* are  
269 two important properties that affect the ability of wavelet bases to approximate  
270 a particular class of functions with few non-zero wavelet coefficients [18]. In  
271 order to reconstruct likelihoods accurately using from a small number of coef-  
272 ficients, we wish wavelet functions to have large number of vanishing moments  
273 and small size of support. Coiflets [19] are a wavelet basis with large number of  
274 vanishing moments and Symmlets [20] are a wavelet basis that have minimum  
275 size of support. The performance of these domains has been analyzed in the  
276 context of our experiments and a proper domain has been selected accordingly  
277 as described in Section 4.2.

## 278 4. Experimental Results

### 279 4.1. Setup

280 In the experiments, we have simulated the VSN environment by using the in-  
281 door multi-camera dataset in [7]. This dataset includes four people sequentially  
282 entering a room and walking around. The sequence was shot by four synchro-  
283 nized cameras in a  $50 m^2$  room. The cameras were located at each corner of the  
284 room. In this sequence, the area of interest was of size  $5.5 m \times 5.5 m \simeq 30 m^2$   
285 and discretized into  $G = 56 \times 56 = 3136$  locations, corresponding to a regular  
286 grid with a  $10cm$  resolution. For the correspondence between camera views and  
287 the top view, the homography matrices provided with the dataset are used. The



Figure 3: A sample set of images from the indoor multi-camera dataset [7].

288 size of the images are  $360 \times 288$  pixels and the frame rate for all of the cameras  
 289 is 25 fps. The sequence is approximately 2.5 minutes ( $\simeq 3,800$  frames) long.

290

291 Starting from the frames around the 2,000th, we have observed failures in the  
 292 original method [7] on preserving identities. For this reason, we have used the  
 293 sequence consisting of the first 2,000 frames for testing. A sample set of images  
 294 is shown in Figure 3.

#### 295 4.2. Comparison of Domains

296 As discussed in Section 3.2, it is very important to select a domain (matrix  
 297  $A$  in Eq. (6)) that can compress the likelihood functions effectively. To select a  
 298 proper domain, we have performed a comparison between DCT, Haar, Symmlet,  
 299 and Coiflet domains and examined the errors in reconstructing the likelihoods  
 300 using various number of coefficients. For the Symmlet domain, the size of sup-  
 301 port is set to 8 and for the Coiflet domain, the number of vanishing moments  
 302 is set to 10. In the comparison, we have used 20 different likelihood functions  
 303 obtained from the tracker in [7]. We have also analyzed the effect of block size  
 304 by choosing two different block sizes:  $8 \times 8$  and  $4 \times 4$ . After we transform each  
 305 block to a domain, we have reconstructed the blocks by using only 1, 2, 3, 4, 5,  
 306 and 10 most significant coefficient(s). In total, for a block size of  $8 \times 8$ , taking  
 307 the most significant 2 coefficients results in 98 coefficients overall. According  
 308 to the structure of the likelihood functions, the elements in a block may all be  
 309 zero. For such a block all the coefficients will be zero, thereby we do not need to  
 310 take coefficients. Thus, we may end up with even smaller number of coefficients.

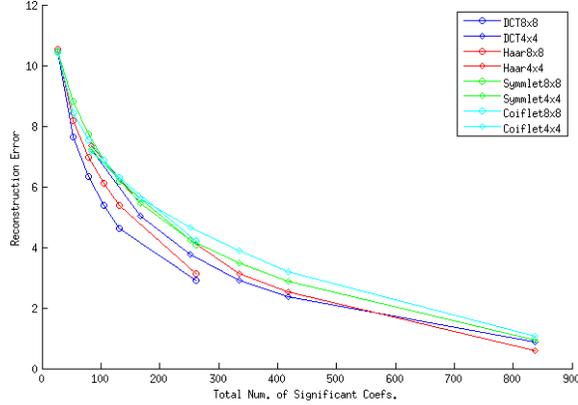


Figure 4: The average reconstruction errors of DCT, Haar, Symmlet, and Coiflet domain for block sizes of  $8 \times 8$  and  $4 \times 4$  using 1, 2, 3, 4, 5 and 10 most significant coefficient(s) per block.

311

312 Figure 4 shows the average of reconstruction errors of each domain for differ-  
 313 ent block sizes. As explained above, the total number of significant coefficients  
 314 used for reconstruction may change depending on the structure of likelihoods.  
 315 For this reason, the x-axis in Figure 4 are the average of number of coefficients  
 316 obtained by taking the 1, 2, 3, 4, 5 and 10 most significant coefficient(s) per  
 317 block. We can see that using DCT with a block size of  $8 \times 8$  outperforms other  
 318 domains. Following this observation, in our tracking experiments, this setting  
 319 has been used.

### 320 4.3. Tracking Results

321 In this subsection, we present the performance of our method used for multi-  
 322 view multi-person tracking. In the experiments, we have compared our method  
 323 with the traditional centralized approach of compressing raw images. In this  
 324 centralized approach, after the raw images are acquired by the cameras, similar  
 325 to JPEG compression, each color channel in the images are compressed and  
 326 sent to the central node. In the central node, features are extracted from the  
 327 reconstructed images and tracking is performed using the method in [7]. For

328 both our method and the centralized approach we have used DCT domain with  
329 a block size of  $8 \times 8$  and took only the 1, 2, 3, 4, 5, 10, and 25 most significant  
330 coefficient(s). Consequently, in our method with the likelihoods of  $56 \times 56$  size,  
331 at each camera in total we end up with at most 49, 98, 147, 196, 245, 490  
332 and 1225 coefficients per person. Since there are four individuals in the scene  
333 at maximum, each camera sends at most 196, 392, 588, 784, 980, 1960 and  
334 4900 coefficients. As mentioned in the previous section, these are the maximum  
335 number of coefficients, since there may be some all-zero blocks. To make a fair  
336 comparison, in the centralized approach we compress the images with  $360 \times 288$   
337 size and 3 color channels. Hence, at each camera we end up with 4860, 9720,  
338 14580, 19440, 24300, 48600 and 121500 coefficients.

339

340 A groundtruth for this sequence is obtained by manually marking the peo-  
341 ple on ground plane, in intervals of 25 frames. Tracking errors are evaluated  
342 via Euclidean distance between the tracking and manual marking results (in  
343 intervals of 25 frames). Figure 5 presents the average of tracking errors over all  
344 people versus the total number of significant coefficients used in communication  
345 for the centralized approach and for our method. Since the total number of sig-  
346 nificant coefficients sent by a camera in our method may change depending on  
347 the structure of likelihood functions and the number of people at that moment,  
348 the maximum is shown in Figure 5. It can be clearly seen that the centralized  
349 approach is not capable of decreasing the communication without affecting the  
350 tracking performance. It needs at least 121500 significant coefficients in total to  
351 achieve an error of around 1 pixel in the grid on average. On the other hand,  
352 our method, down to using 3 significant coefficients per block, achieves an error  
353 of around 1 pixel in the grid on average. In our experiments, this led to sending  
354 at most 408 coefficients for four people. Taking less than 3 coefficients per block  
355 affects the performance of the tracker and produces an error of 11.5 pixels in  
356 the grid on average. But in overall, our method significantly outperforms the  
357 centralized approach.

358

359 The tracking errors for each person and the tracking results, obtained by the  
360 centralized approach using 48600 coefficients in total, are given in Figure 6-  
361 a and Figure 6-b, respectively. It can be seen that although the centralized  
362 approach can track the first and the second individuals very well, there is an  
363 identity association problem for the third and fourth individuals. In Figure 7-a  
364 and Figure 7-b, we present the tracking errors for each person and the tracking  
365 results obtained with our method using 3 coefficients per block, respectively.  
366 Clearly, we can see that all people in the scene can be tracked very well by our  
367 method. The reason of the peak error value in the third person is because the  
368 tracking starts a few frames after the third person enters the room. For this  
369 reason, there is a big error at the time third person enters the room. When the  
370 number of coefficients taken per block are less than 3, we also observe identity  
371 problems. But by selecting the number of coefficients per block greater than or  
372 equal to 3, we can track all the people in the scene accurately. The centralized  
373 approach, in total, requires at least more than two orders of magnitude coeffi-  
374 cients to achieve this level of accuracy.

375

376 In the light of the results we obtained, for the same tracking performance,  
377 our framework saves 99.6% of the bandwidth compared to the centralized ap-  
378 proach. Our framework is also advantageous over an ordinary decentralized  
379 approach that directly sends likelihood functions to the fusion node. In such  
380 an approach, we send each data point in the likelihood function, resulting a  
381 need of sending 12544 values for tracking four people. The performance of this  
382 approach is also given in Figure 5. For the same level of tracking accuracy, our  
383 framework achieves saving 96.75% compared to the decentralized approach.

## 384 5. Conclusion

385 Visual sensor networks constitute a new paradigm that merges two well-  
386 known topics: computer vision and sensor networks. Consequently, it poses  
387 unique and challenging problems that do not exist either in computer vision or

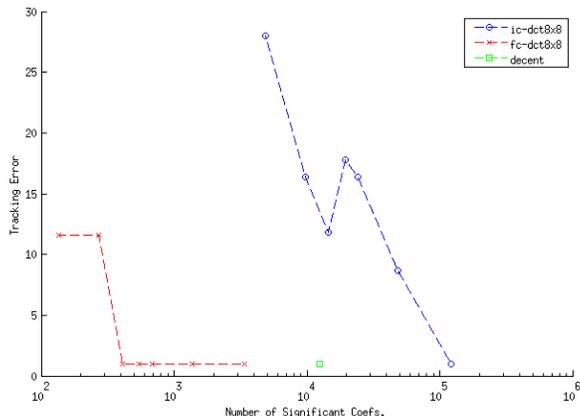
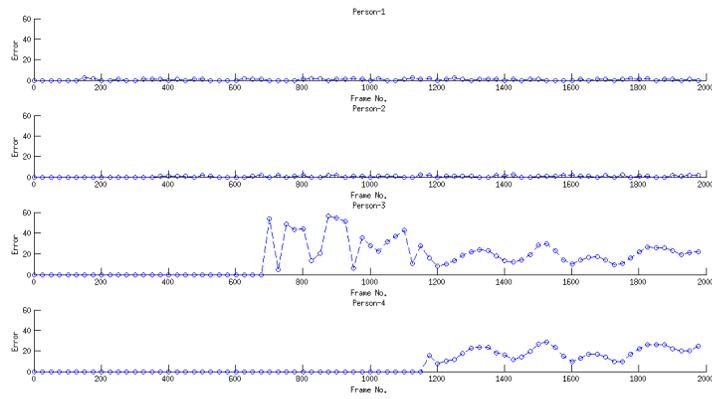
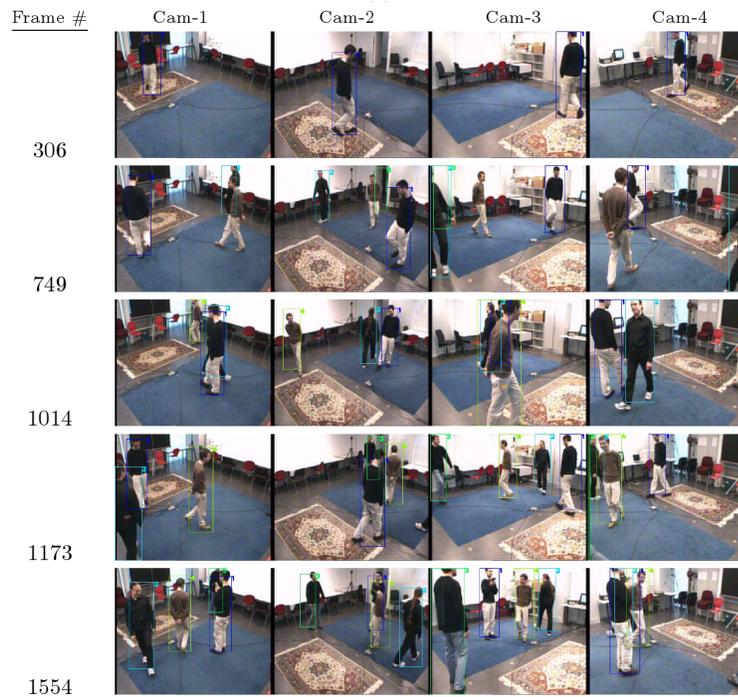


Figure 5: The average tracking errors of the centralized approach (“ic-dct8x8”), our framework (“fc-dct8x8”) both using DCT with  $8 \times 8$  blocks and a decentralized method (“decent”) that directly sends likelihood functions versus the total number of significant coefficients used in reconstruction.

388 in sensor networks. This paper presents a novel method that can be used in  
389 VSNs for multi-camera person tracking applications. In our framework, track-  
390 ing is performed in a decentralized way: each camera extracts useful features  
391 from the images it has observed and sends them to a fusion node which collects  
392 the multi-view image features and performs tracking. In tracking, extracting  
393 features usually results a likelihood function. Instead of sending the likelihood  
394 functions itself to the fusion node, we compress the likelihoods by first splitting  
395 them into blocks, and then transforming each block to a proper domain and tak-  
396 ing only the most significant coefficients in this representation. By sending the  
397 most significant coefficients to the fusion node, we decrease the communication  
398 in the network. At the fusion node, the likelihood functions are reconstructed  
399 back and tracking is performed. The idea of performing goal-directed compression  
400 in a VSN is the main contribution of this work. Rather than focusing on  
401 low-level communication without regard to the final inference goal, we propose a  
402 different compressing scheme that is better matched to the final inference goal,  
403 which, in the context of this paper, is tracking.

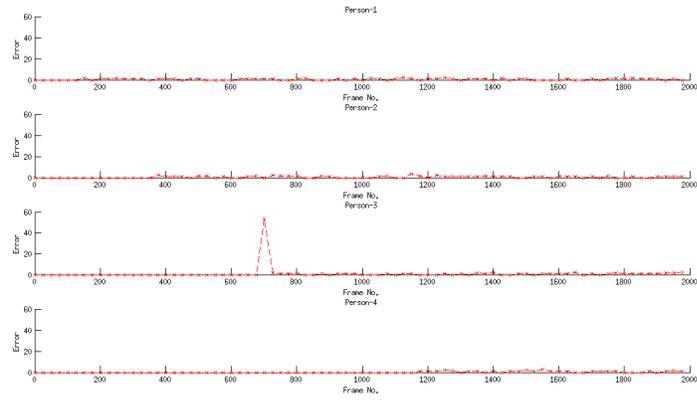


(a)

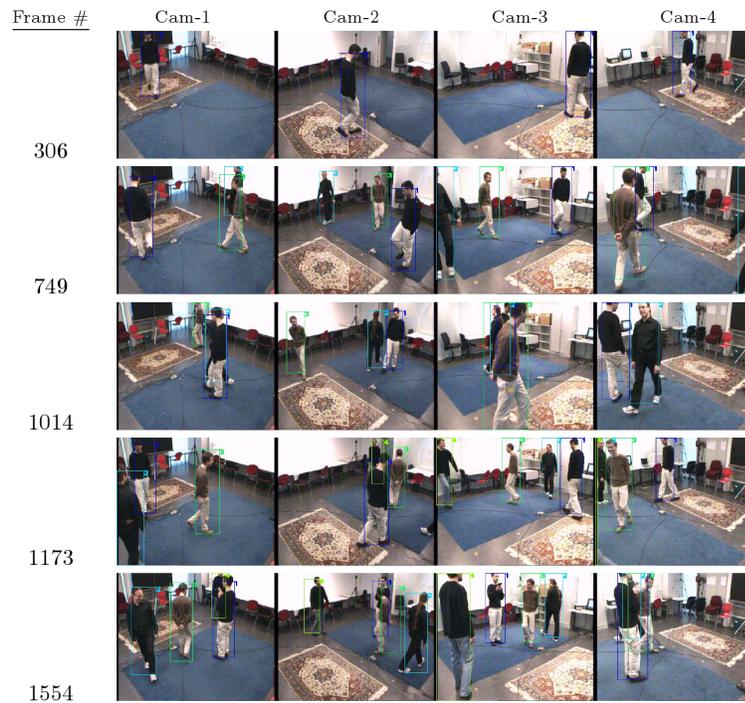


(b)

Figure 6: (a) The tracking errors for each person and (b) tracking results obtained by the centralized approach using 48600 coefficients in total used in communication.



(a)



(b)

Figure 7: (a) The tracking errors for each person and (b) tracking results obtained by our framework using 3 coefficients per block used in communication.

404

405 This framework fits well to the needs of the VSN environment in two aspects: i)  
406 the processing capabilities of cameras in the network are utilized by extracting  
407 image features at the camera-level, ii) using only the most significant coeffi-  
408 cients in network communication saves energy and bandwidth resources. We  
409 have achieved a goal-directed compression scheme for the tracking problem in  
410 VSNs by performing local processing at the nodes and compressing the resulting  
411 likelihood functions which are related to the tracking goal, rather than compress-  
412 ing raw images. To the best of our knowledge, this method is the first method  
413 that compresses likelihood functions and applies this idea for VSNs. Another  
414 advantage of this framework is that it does not require the use of a specific track-  
415 ing method. Without making significant changes on existing tracking methods  
416 (e.g., using simpler features, etc.), which may degrade the performance, such  
417 methods can be used within our framework in VSN environments. In the light  
418 of the experimental results, we can say that our feature compression approach  
419 can be used together with any robust probabilistic tracker in the VSN context.

420

421 We believe that trying different dictionaries that are better matched to the  
422 structure of likelihood functions, thereby, leading to further reductions in the  
423 communication load, can be a possible direction for future work. In addition,  
424 an interesting future work direction can be the implementation of our method  
425 in a real VSN setup.

#### 426 **Acknowledgements**

427 This work was partially supported by a Turkish Academy of Sciences Distin-  
428 guished Young Scientist Award and by a graduate scholarship from the Scientific  
429 and Technological Research Council of Turkey.

430 **References**

- 431 [1] P. V. Pahalawatta, A. K. Katsaggelos, Optimal sensor selection for video-  
432 based target tracking in a wireless sensor network, in: in Proc. International  
433 Conference on Image Processing (ICIP '04, 2004, pp. 3073–3076.
- 434 [2] S. Fleck, F. Busch, W. Straß er, Adaptive probabilistic track-  
435 ing embedded in smart cameras for distributed surveillance in a  
436 3d model, EURASIP J. Embedded Syst. 2007 (1) (2007) 24–24.  
437 doi:http://dx.doi.org/10.1155/2007/29858.
- 438 [3] E. Oto, F. Lau, H. Aghajan, Color-based multiple agent tracking for wire-  
439 less image sensor networks, in: ACIVS06, 2006, pp. 299–310.
- 440 [4] B. Song, A. Roy-Chowdhury, Robust tracking in a camera net-  
441 work: A multi-objective optimization framework, Selected Topics  
442 in Signal Processing, IEEE Journal of 2 (4) (2008) 582 –596.  
443 doi:10.1109/JSTSP.2008.925992.
- 444 [5] H. Medeiros, J. Park, A. Kak, Distributed object tracking using a  
445 cluster-based kalman filter in wireless camera networks, Selected Top-  
446 ics in Signal Processing, IEEE Journal of 2 (4) (2008) 448 –463.  
447 doi:10.1109/JSTSP.2008.2001310.
- 448 [6] J. Yoder, H. Medeiros, J. Park, A. Kak, Cluster-based distributed face  
449 tracking in camera networks, Image Processing, IEEE Transactions on  
450 19 (10) (2010) 2551 –2563. doi:10.1109/TIP.2010.2049179.
- 451 [7] F. Fleuret, J. Berclaz, R. Lengagne, P. Fua, Multicamera peo-  
452 ple tracking with a probabilistic occupancy map, PAMI 30 (2).  
453 doi:10.1109/TPAMI.2007.1174.
- 454 [8] M. Taj, A. Cavallaro, Distributed and decentralized multicamera  
455 tracking, Signal Processing Magazine, IEEE 28 (3) (2011) 46 –58.  
456 doi:10.1109/MSP.2011.940281.

- 457 [9] J. Yao, J.-M. Odobez, Multi-camera multi-person 3d space tracking with  
458 mcmc in surveillance scenarios, in: ECCV workshop on Multi Camera and  
459 Multi-modal Sensor Fusion Algorithms and Applications, 2008.
- 460 [10] A. Gupta, A. Mittal, L. Davis, Constraint integration for effi-  
461 cient multiview pose estimation with self-occlusions, PAMI 30 (3).  
462 doi:10.1109/TPAMI.2007.1173.
- 463 [11] M. Hofmann, D. Gavrilu, Multi-view 3d human pose estimation combin-  
464 ing single-frame recovery, temporal integration and model adaptation, in:  
465 CVPR, 2009. doi:10.1109/CVPR.2009.5206508.
- 466 [12] D. Comaniciu, V. Ramesh, P. Meer, Kernel-based object tracking, Pattern  
467 Analysis and Machine Intelligence, IEEE Transactions on 25 (5) (2003) 564  
468 – 577. doi:10.1109/TPAMI.2003.1195991.
- 469 [13] T.-L. Liu, H.-T. Chen, Real-time tracking using trust-region methods, Pat-  
470 tern Analysis and Machine Intelligence, IEEE Transactions on 26 (3) (2004)  
471 397 –402. doi:10.1109/TPAMI.2004.1262335.
- 472 [14] P. Perez, C. Hue, J. Vermaak, M. Gangnet, Color-based probabilistic  
473 tracking, in: Heyden, A and Sparr, G and Nielsen, M and Johansen, P  
474 (Ed.), COMPUTER VISION - ECCV 2002, PT 1, Vol. 2350 of LECTURE  
475 NOTES IN COMPUTER SCIENCE, IT Univ Copenhagen; Univ Copen-  
476 hagen; Lund Univ, 2002, pp. 661–675, 7th European Conference on Com-  
477 puter Vision (ECCV 2002), COPENHAGEN, DENMARK, MAY 28-31,  
478 2002.
- 479 [15] G. Wallace, The jpeg still picture compression standard, Con-  
480 sumer Electronics, IEEE Transactions on 38 (1) (1992) xviii –xxxiv.  
481 doi:10.1109/30.125072.
- 482 [16] M. Antonini, M. Barlaud, P. Mathieu, I. Daubechies, Image coding using  
483 wavelet transform, Image Processing, IEEE Transactions on 1 (2) (1992)  
484 205 –220. doi:10.1109/83.136597.

- 485 [17] L. Winger, A. Venetsanopoulos, Biorthogonal modified coiflet filters for  
486 image compression, in: Acoustics, Speech and Signal Processing, 1998.  
487 Proceedings of the 1998 IEEE International Conference on, Vol. 5, 1998,  
488 pp. 2681 –2684 vol.5. doi:10.1109/ICASSP.1998.678075.
- 489 [18] S. Mallat, A Wavelet Tour of Signal Processing, Second Edition (Wavelet  
490 Analysis & Its Applications), 2nd Edition, Academic Press, 1999.
- 491 [19] I. Daubechies, Orthonormal bases of compactly supported wavelets, Com-  
492 munications on Pure and Applied Mathematics.
- 493 [20] I. Daubechies, Ten lectures on wavelets, 1st Edition, Society for Industrial  
494 and Applied Mathematics, 1992.